

DOCUMENT RESUME

ED 164 301

SE 025 456

AUTHOR Beck, A.; And Others
TITLE Calculus, Part 1, Student's Text, Unit No. 66.
Revised Edition.
INSTITUTION Stanford Univ., Calif. School Mathematics Study
Group.
SPONS AGENCY National Science Foundation, Washington, D.C.
PUB DATE 66
NOTE 373p.; For related documents, see SE 025 457-459;
Contains occasional light type

EDRS PRICE MF-\$0.83 HC-\$19.41 Plus Postage.
DESCRIPTORS *Calculus; *Curriculum; *Instructional Materials;
*Mathematical Applications; Mathematics Education;
Secondary Education; *Secondary School Mathematics;
*Textbooks
IDENTIFIERS *School Mathematics Study Group

ABSTRACT

This is part one of a three-part SMSG calculus text for high school students. One of the goals of the text is to present calculus as a mathematical discipline as well as presenting its practical uses. The authors emphasize the importance of being able to interpret the concepts and theory in terms of models to which they apply. The text demonstrates the origins of the ideas of the calculus in practical problems; attempts to express these ideas precisely and develop them logically; and finally, returns to the problems and applies the theorems resulting from that development. Chapter topics include: (1) Introduction; (2) The Idea of Derivative; (3) Limits and Continuity; (4) Differentiation; and (5) Applications of the Derivative. (MP)

* Reproductions supplied by EDRS are the best that can be made *
* from the original document. *

Calculus

Part 1 Student's Text

REVISED EDITION

The following is a list of all those who participated in the preparation of this volume:

A. Beck	Olney High School, Philadelphia, Pa.
A. A. Blank	New York University, New York, N.Y.
F. L. Elder	West Hempstead Jr., Sr. High School, N.Y.
C. E. Kerf	Dickinson College, Carlisle, Pa.
M. S. Klamkin	Ford Scientific Laboratory, Dearborn, Mich.
I. I. Kolodner	Carnegie Institute of Technology, Pittsburgh, Pa.
M. D. Kruskal	Princeton University, Princeton, N.J.
C. W. Leeds, III	Berkshire School, Sheffield, Mass.
M. A. Linton, Jr.	William Penn Charter School, Philadelphia, Pa.
H. M. Marston	Douglass College, New Brunswick, N.J.
I. Marx	Purdue University, Lafayette, Ind.
R. Pollack	New York University, New York, N.Y.
T. L. Reynolds	College of William and Mary, Williamsburg, Va.
R. L. Starkey	Cubberley High School, Palo Alto, Calif.
V. Twersky	Sylvania Electronics Defense Labs., Mt. View, Calif.
H. Weitzner	New York University, New York, N.Y.

New Haven and London, Yale University Press, 1966

Financial support for the School Mathematics Study Group has been provided by the National Science Foundation.

Permission to make verbatim use of material in this book must be secured from the Director of SMSG. Such permission will be granted except in unusual circumstances. Publications incorporating SMSG materials must include both an acknowledgment of the SMSG copyright (Yale University or Stanford University, as the case may be) and a disclaimer of SMSG endorsement. Exclusive license will not be granted save in exceptional circumstances, and then only by specific action of the Advisory Board of SMSG.

© 1965 by The Board of Trustees
of the Leland Stanford Junior University.
All rights reserved.
Printed in the United States of America.

PREFACE

This calculus is offered to you as a book, not merely as a manual of instructions. There is no intention of presenting the subject as a loosely strung-together sequence of topics composed only as a crib to help you recite the "right" answers to traditional examination questions.

The authors have tried to convey their own attitude about the subject as it is regarded maturely; a development based upon a very small number of major themes, with many lesser themes all tightly interwoven with the major ones. This sense of plot of total structure is worth your effort to comprehend. When the calculus is seen in this perspective, the multitudinous details no longer confuse but take their proper place in relation to the grand structure.

Although calculus is an old part of the body of human thought it remains alive and useful. To us, the authors, who have had to think hard about what we wished to say about the calculus (and even whether it was worth adding yet another book to the dense population), it has been a high intellectual adventure. Even though the subject is old and well-explored, each of us gained new insights in re-creating it for himself. We envy you, to whom the subject is entirely new; you will have the pleasure of re-creating it entire. We hope that you will find the task of re-creation as we did, always absorbing, sometimes amusing, and often exciting.

FOREWORD

Use of the Text

The sections, definitions, theorems, examples, figures, and exercises in the text are labeled according to the chapter. For example, Section 6-4 is the fourth section in Chapter 6, and Example 6-4c is the third example in Section 6-4.

If only one corollary follows a theorem or lemma the corollary is not numbered. If more than one follows, then the corollaries are numbered in order without numbered reference to the theorem; thus, Theorem 5-2b is followed by Corollaries 1 and 2.

Sections and exercises marked with the symbol Λ (for a Himalayan peak) will generally present an extra challenge. You will want to work through this material "because it is there."

It is expected that a handbook of mathematical tables will be readily available and, hence, no extensive reproduction of tables and formulas is provided here.

A general summary of prerequisite material is given in appendices.

Some of the symbols used in the text usually have a restricted denotation (e.g., the early letters of the alphabet are generally reserved for constants, n for an integer). We shall use the Greek alphabet (Appendix O) as well as the English.

Appendix 1 treats the general properties of real numbers. A thorough understanding of those properties discussed in Sections A1-2, A1-3, and A1-4 is essential to the work on limits. For most readers, the material in Section A1-4, "Intervals, Neighborhoods," is new and will require special attention. Section A1-5 treats the completeness of the real number system and is essential for reading in appendices which amplify the text (App. 4, et seq.).

Appendix 2 summarizes basic ideas involving functions and the properties of special functions. The text presupposes familiarity with functional terminology and notation. In particular, Example 2-5c and Section 4-4 assume knowledge of circular (trigonometric) functions.

Appendix 3 gives an account of proof by mathematical induction, a method which will be useful in the solutions of many exercises (e.g., Exercises 3-4, No. 1).

Sum notation (A3-2) is first used in Section 6-1. A working knowledge of the symbolism is essential for facility in the work on integration.

The appendices are also used to supply information which is not ordinarily given in a first calculus course. For the sake of more complete understanding of the subject, you may wish to read some of this material.

TABLE OF CONTENTS

PREFACE	1
FOREWORD	ii
Chapter 1. INTRODUCTION	1
1-1. "Best-Value" Problems. The Derivative	2
1-2. Area. The Integral	11
1-3. The Scope of the Calculus	21
Chapter 2. THE IDEA OF DERIVATIVE	27
2-1. Introduction	27
2-2. Slope	30
2-3. General Quadratic Function	37
2-4. Velocity	42
2-5. Derivative of a Function at a Point	49
Chapter 3. LIMITS AND CONTINUITY	55
3-1. Introduction	55
3-2. Definition of Limit of a Function	58
3-3. Epsilonic Technique	67
3-4. Limit Theorems	78
3-5. The Idea of Continuity	91
3-6. Properties of Functions Continuous at a Point	99
3-7. Properties of Functions Continuous on an Interval	108
Chapter 4. DIFFERENTIATION	117
4-1. Introduction	117
4-2. Rational Combinations	120
4-3. Inverse Functions, Fractional Powers	131
4-4. Circular Functions	137
4-5. Inverse Circular Functions	143
4-6. Compositions. Chain Rule	149
4-7. Notation	154
4-8. Implicitly Defined Functions	161
Miscellaneous Exercises	167

Chapter 5. APPLICATIONS OF THE DERIVATIVE	169
5-1. Introduction	169
5-2. The Derivative at an Extremum	173
5-3. The Law of the Mean	186
5-4. Applications of the Law of the Mean	196
5-5. Applications of the Second Derivative	205
5-6. Constrained Extreme Value Problems	213
5-7. Tangent and Normal Lines	223
5-8. Sketching of Graphs	229
Miscellaneous Exercises	240
Appendix O. THE GREEK ALPHABET	244
Appendix 1. THE REAL NUMBERS	245
A1-1. Algebraic Properties of the Real Numbers	245
A1-2. Order Relations on the Real Numbers (Inequality)	249
A1-3. Absolute Value and Inequality	254
A1-4. Intervals, Neighborhoods	259
A1-5. Completeness of the Real Number System. Separation Axiom	263
Appendix 2. FUNCTIONS AND THEIR REPRESENTATIONS	269
A2-1. Functions	269
A2-2. Composite Functions	285
A2-3. Inverse Functions	290
A2-4. Monotone Functions	298
A2-5. The Circular Functions	303
A2-6. Polar Coordinates	308
Appendix 3. MATHEMATICAL INDUCTION	319
A3-1. The Principle of Mathematical Induction	319
A3-2. Sums and Sum Notation	333
Appendix 4. FUNCTIONS CONTINUOUS ON AN INTERVAL	345
A4-1. The Extreme Value Theorem	345
A4-2. The Intermediate Value Theorem	350
A4-3. A Nowhere Differentiable Continuous Function	352
Appendix 5. IMPLICITLY DEFINED FUNCTIONS AND THEIR DERIVATIVES	359

Chapter 1

INTRODUCTION

The calculus is a powerful and flexible instrument for obtaining useful solutions to an astonishing variety of problems in science, technology, and industry. It is also a mathematical discipline in which theorems are deduced from carefully stated postulates and definitions and for which faultless logic is primary. These are complementary aspects of the subject. To be capable of making the most efficient practical use of the calculus it is important to understand the reasoning upon which its techniques are based. To understand why the concepts and theory of the calculus are significant, even to care about developing the theory in the first place, it is important to be able to interpret the concepts and theory in terms of models, whether geometrical, physical or other, to which they apply. In this text we shall find the origins of the ideas of the calculus in practical problems; we shall attempt to express these ideas precisely so that we may reason about them logically; finally, we shall return to problems and apply the theorems resulting from our reasoning.

The two basic ideas of the elementary calculus are "derivative" and "integral". It is easy to appreciate these ideas intuitively and know why they are useful before formulating them precisely. Here we shall consider these ideas as they arise in the solution of specific problems.

1-1. "Best-Value" Problems. The Derivative.

It is in the nature of the human enterprise to try to get the best of everything: a manufacturer seeks the smallest unit cost for his product and the highest possible price; a student tries to complete his homework assignment in the shortest possible time; a demagogue expounds the political philosophy which he believes will garner the greatest number of votes. It is seldom clear what must be done to get the best value. We may even fail to make the observation that we have a problem of this kind. Immediately below the surface of daily life there are countless "best-value" problems which are often overlooked because each our consciousness cloaked in a wealth of detail which must be stricken away before the central problem is perceived. Not all of these problems can be solved by the calculus alone, but many of them can, and it is just such a problem which we now consider.

In the solution of this problem, we shall go far beyond the casual accident of humdrum experience from which it came. The solution introduces some of the deepest ideas in the history of human thought and in understanding the solution you will gain an intuitive appreciation of these ideas. Furthermore, the systematic approach by which the problem is solved will be appreciated not only for the particular solution of this very specific problem, but for the possibility of using the same method for the solutions of an extensive, important class of "best-value" problems. Eventually (Chapter 5), we shall see how this systematic attack can be sharpened to give an extremely simple and direct method of calculation for obtaining best values.

This is a problem from my own experience. I am fond of books to the point of avarice and in time have acquired a substantial number. Several years ago, when I had to move my household from one city to another I was appalled at the cost of moving the enormous dead weight of my books by interstate van. Upon investigating a number of alternatives, I found that there is a special book rate for parcel post and that this offered by far the cheapest method of shipping books. The post office restricts the size of parcels. To keep the effort of packing to a minimum, I bought out the largest possible cartons complying with the post office requirement at that time.

The post office regulations for parcel post state, "Parcels mailed at a first-class post office in the United States for delivery at ... any ... first-class post office ... must not exceed 72 inches in length and girth combined." To apply this regulation we must know that the length of a rectangular carton is that of its longest edge and the girth is the perimeter of the cross section perpendicular to that edge. Our problem is to determine the dimensions of the carton of largest size complying with the post office regulation.

For simplicity we consider only cartons with a square cross section.* If, in particular, the box were a cube and e the length of an edge in inches then the girth would be $4e$, and for the cube of maximum size permitted by postal regulations

$$e + 4e = 72$$

or
$$e = \frac{72}{5}.$$

Does a cube give us the largest carton that can be sent by parcel post?

In order to answer this question we let x denote the length in inches of a side of the square cross section and y the length in inches of the longest side of the carton, so $y \geq x$. Taking x as large as possible for a given y we require that

$$(1) \quad y + 4x = 72, \quad (y \geq x).$$

Under these conditions we attempt to maximize the volume of the carton,

$$V = x^2 y.$$

Setting $y = 72 - 4x$ in the expression for V we obtain

$$(2) \quad V = x^2(72 - 4x),$$

a formula valid so long as $y \geq x \geq 0$, that is, from (1), so long as $0 \leq x \leq \frac{72}{5}$. These conditions define a function $f: x \rightarrow V$.

Our problem is not so much to determine the largest value V_{\max} in the range of the function, although that information may also be useful, but to find a value of a in the domain of f for which $f(a) = V_{\max}$ --the maximum value of the volume is less relevant for our purpose than the dimensions of the box having that capacity.

In order to estimate the location of a maximum value, we sketch a graph of the function. The following table of values (extending beyond $x = \frac{72}{5}$ for convenience in sketching) gives us the coordinates of a few easily calculated points of the graph of f :

*To see that a square cross section is best for a given girth g , we set $s = g/4$ and obtain expressions for adjacent sides of a rectangular section in the form $s - c$ and $s + c$. For a given girth, hence, for a given s , the area $s^2 - c^2$ of a cross section is clearly largest when $c = 0$. It follows that for a given girth the volume is largest when the cross section is square.

x	0	1	5	10	15	18
v	0	68	1300	3200	2700	0

We plot the corresponding points of the graph and sketch a smooth curve, joining them. (unbroken curve in Figure 1-1a).

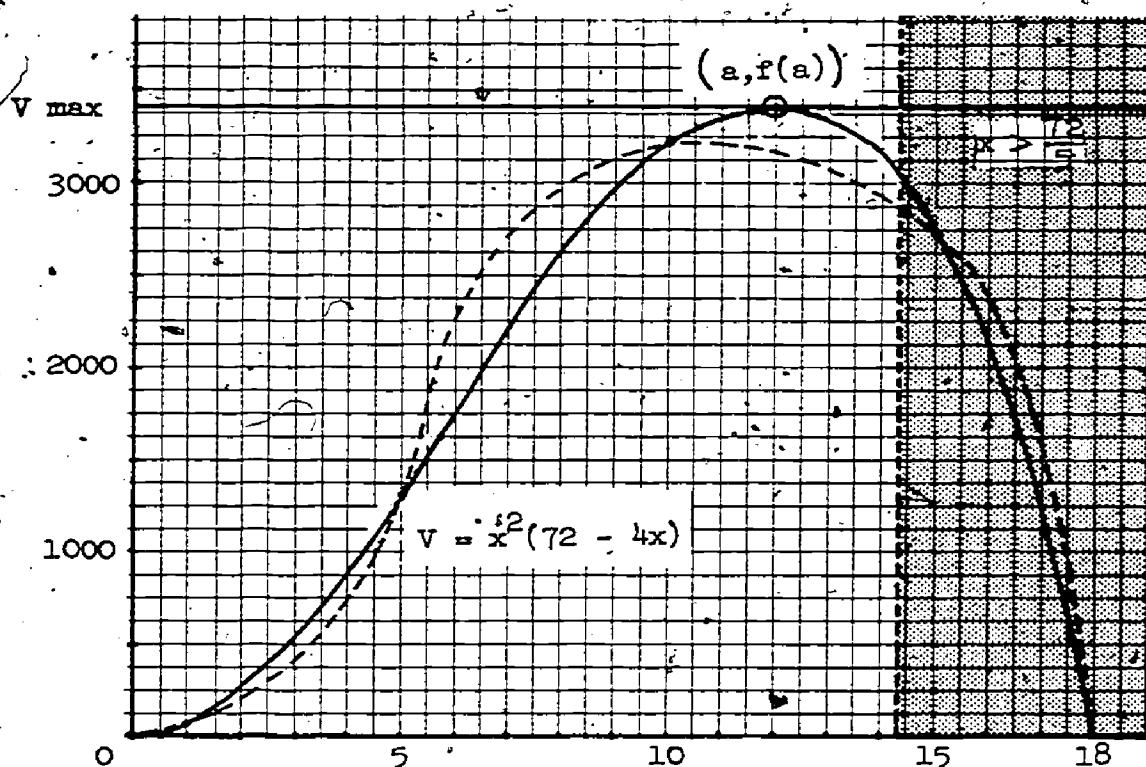


Figure 1-1a

The curve does suggest the approximate location of a peak of the graph and the table does give some precise information, such as

$$V_{\max} \geq f(10) = 3200.$$

No matter how much information we get this way we shall always be somewhat dissatisfied. In the first place, we have exact information about the function only at a few calculated points so that even if we stumbled upon the maximum we might not be aware of it. In the second place, the idea of drawing a smooth curve through the calculated points has its limitations. For example, without further calculation we could not be sure that the unbroken curve in Figure 1-1a more reasonably represents the function than the dashed one, and, furthermore, we cannot eliminate this kind of ambiguity completely by calculating more points. One of our objectives is to devise systematic methods for resolving these difficulties.

Thinking of the problem in visual geometrical terms we see that the condition for a maximum, $f(a) = V_{\max}$, means that the graph of f cannot cross over the horizontal line through $(a, f(a))$ to a higher point. The direction of the graph at $(a, f(a))$ must therefore also be horizontal, for if the graph met the line at an angle, the two would have to cross.* We conclude that to locate a peak of the graph, we seek it among the points where the direction of the graph is horizontal.

In order to make some general use of this geometrical idea we express it numerically so that it may serve as a basis for computation. The direction of the graph at $(a, f(a))$ is determined by the angle the graph makes with the line $y = f(a)$ parallel to the x -axis. For our purposes the direction is most conveniently measured by the tangent of this angle, the slope of the graph. We represent the direction of the graph numerically in terms of slope and reformulate our idea: at a peak of the graph the slope is zero.

We introduce the function f' where $f'(x)$ is the slope of the graph of f at the point $(x, f(x))$. The function f' is called the derivative of f (meaning "derived from f "), and the slope $f'(x)$ of the graph of f at $(x, f(x))$ is called the derivative of f at x . When there is a peak of the graph of f at $(a, f(a))$, we have $f'(a) = 0$; to locate a peak, then, we look among the zeros of f' .

Turning back to our original problem we find that thus far we have only replaced it by new problems. In particular, we have not clearly defined the slope of the graph of f at a given point, the derivative at the point. Furthermore, even if the derivative at the point is defined, there remains the problem of describing the function f' in such a way that we can find its zeros in order to locate a peak of the graph of f .

By now you may feel that we are very far from the point of beginning and that you would like to know what we have accomplished. What we have done is this: we have replaced a problem about which we know very little with a problem about which we know a great deal: to locate a peak of one function we look among the zeros of another function (the derivative). It may seem to you that the line of approach is indirect and it is still not clear that it will be fruitful. We promise that it will be fruitful. You should not think that the discovery of such an avenue of investigation requires superhuman powers. Whenever you become unduly impressed by the ingenuity and power of mathematical methods, reflect that an investigator will try

* Here we are making an outright assumption, that the graph has a definite direction at each point. In particular, the graph may not have a sharp corner at $(a, f(a))$.

not one but many approaches. To his admiring audience he will present the one idea that worked and never mention the failures that filled his waste basket with reams of paper. In fact, we briefly considered and rejected one idea already, that of finding the maximum value of $f(x)$ by examining a number of its values.

Before we go on to solve our best-value problem, it should be said openly that the method of solution we rejected was a perfectly practical one. Without knowing what you are now learning about such problems, you might have proceeded by calculating values and come very close to the optimum solution.* The point is that problems of this kind arise often, and if we have a great many similar problems it pays to devote some attention to more refined methods of solution. Similarly, if you wished to make just one pin you would be content to do it by hand, but if you wished to produce pins by the million you would put a great deal of effort into designing suitable machinery for the purpose. You will soon learn methods that will make the solution of our present problem appear no more consequential than the production of a single pin in the operation of a pin factory.

To attack the problem of defining the derivative we resort to a standard method of the calculus.

We seek a number. This number will be described by approximations. The set of approximations must be adequate; that is, there must always be an approximation which is closer to the number than any error we may tolerate, no matter how small the specified error tolerance. In the language of the calculus we say the number is the limit of the set of approximations.

To approximate the slope of the graph of f at a point $(a, f(a))$ we consider the arc of the graph between the point $(a, f(a))$ and another point $(x, f(x))$. The statement that the graph of f at $(a, f(a))$ has a certain slope $f'(a)$ will mean now that it is possible to approximate $f'(a)$ closely by the slope of the chord between $(x, f(x))$ and $(a, f(a))$. More precisely, the error in approximating $f'(a)$ by the slope of the chord can be reduced below any given tolerance by taking x close enough to a (Figure 1-1b).

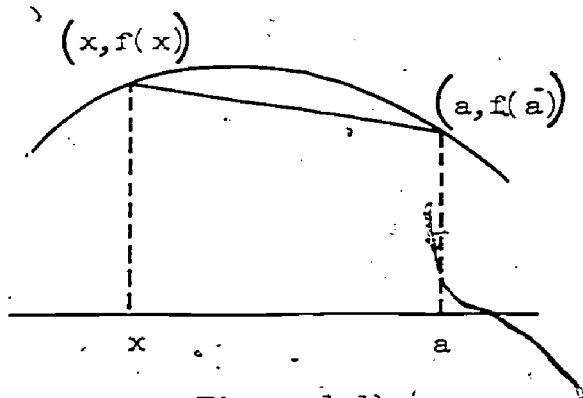


Figure 1-1b

* Since the graph is nearly horizontal in the neighborhood of a peak, the penalty for missing the exact location of the peak can be expected to be quite small. We shall return to this point later in the text.

Now we are ready to attack the box problem directly. For $f(x) = x^2(72 - 4x)$ the slope of the chord between $(x, f(x))$ and $(a, f(a))$ is the ratio

$$\begin{aligned}\frac{f(x) - f(a)}{x - a} &= \frac{x^2(72 - 4x) - a^2(72 - 4a)}{x - a} \\ &= \frac{72(x^2 - a^2) - 4(x^3 - a^3)}{x - a} \\ &= \frac{x - a}{x - a} [72(x + a) - 4(x^2 + ax + a^2)].\end{aligned}$$

When $x = a$ the ratio is algebraically meaningless since the denominator is zero. This is to be expected since the geometrical interpretation of the ratio as the slope of the chord joining two points loses its meaning if $(a, f(a))$ and $(x, f(x))$ represent the same point. For any value of x other than a , we note that we have

$$\frac{x - a}{x - a} = 1,$$

and the other factor

$$72(x + a) - 4(x^2 + ax + a^2)$$

is a polynomial which at $x = a$ has the value $144a - 12a^2$. We shall prove later for a polynomial function $p(x)$ that it is possible to approximate $p(a)$ to within any fixed margin of error by taking x sufficiently close to a . It follows that the slope of the graph at $(a, f(a))$ is

$$f'(a) = 144a - 12a^2.$$

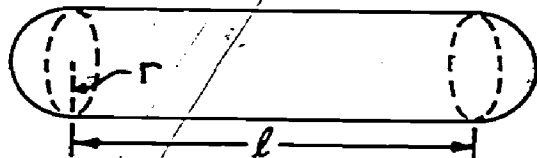
Now we use the criterion that the slope of the graph at a peak is zero. The zeros of f' occur at $a = 0$ and $a = 12$; clearly, $f(0) = 0$ is not the best value. Having eliminated every other possibility, we see that the desired maximum must occur at $x = 12$. In conclusion, the largest carton with length plus girth of 72 inches has square ends with 12-inch sides and a length of 24 inches.

You will have noticed that the actual computation leading to the solution is quite short. Most of the effort and time was spent in explaining the considerations underlying this method of solution. Later we shall see that we may write out $f'(a)$ on sight and the labor of solution will then be almost negligible. Finally, we have produced yet another problem: if we want to find the maximum of $f(x)$, and we know the zeros of f' , then which of these zeros--if any--yields the best value we are seeking? This question we shall leave to be answered later in the text.

Exercises 1-1

This set of exercises is to provide practice in formulating applied problems mathematically.

- Express the area of a semicircle as a function of its perimeter.
- A rectangle is inscribed in a circle of radius R . Express the area of the rectangle as a function of the length of one side.
- A tank has the form of a right circular cylinder capped by hemispherical ends. If its total volume is V , express the length of the cylinder as a function of its radius.



- Determine the edge of a regular tetrahedron inscribed in a sphere as a function of the radius R .

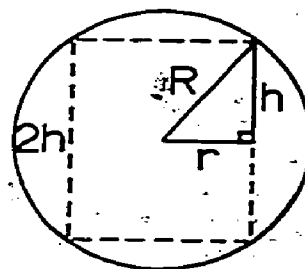
DIRECTIONS: Numbers 5-11. Translate the given information in such a way as to discover a function f which may be maximized or minimized and write an equation defining $f(x)$. (You will be asked to solve the problem later on in the course.)

Illustrative Example:

Find the right circular cylinder of greatest volume that can be inscribed in a sphere of radius R .

Solution:

Consider a plane cross-section of the figure containing the axis of the cylinder and let r and $2h$ denote the radius and height of the cylinder, respectively. The volume of this cylinder is



$$V = 2\pi hr^2$$


where

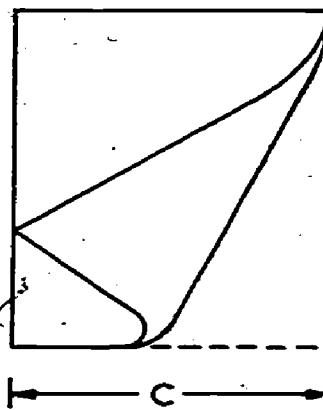
$$r^2 + h^2 = R^2$$

Thus we have to maximize

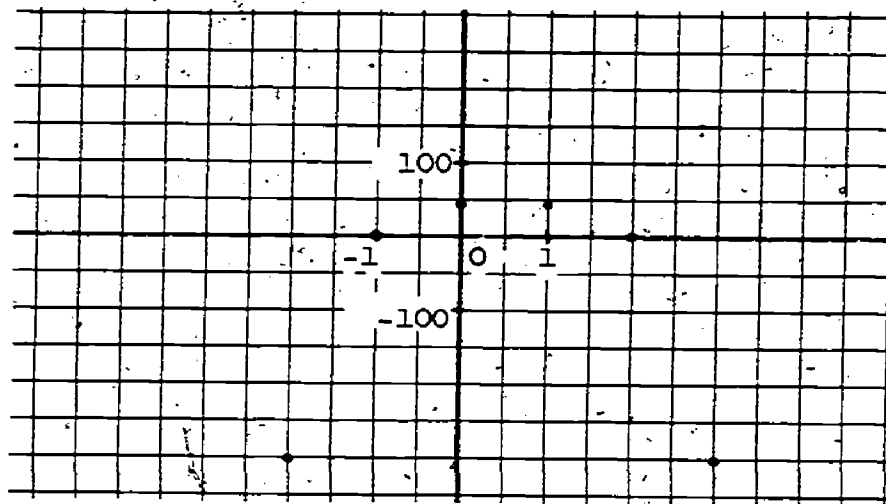
$$V = f(h) = 2\pi h(R^2 - h^2), \quad 0 \leq h \leq R.$$

- What point on the ellipse $x^2 + 4y^2 = 9$ is nearest the point $(1, 6)$?

6. A rectangle has two of its vertices on the x-axis and the other two above the axis on the parabola $y = 6 - x^2$. What are the dimensions of such a rectangle if its area is to be a maximum?
7. A rectangular sheet of galvanized metal is bent to form the sides and bottom of a trough so that the cross section has this shape: . If the metal is 14 inches wide, how deep must the trough be to carry the most water?
8. The tank described in exercise 3 above is to be made in the most economical way. If the material for the hemispherical ends costs twice as much per square foot as that for the cylindrical part, find the most economical dimensions. (That is, find "cost" as a function of either the radius or the length.)
9. Find the right circular cylinder of greatest volume that can be inscribed in a right circular cone of radius r and height h .
10. Find the length of the longest rod which can be carried horizontally around a corner from a corridor 10 ft. wide into one 5 ft. wide.
11. The lower right-hand corner of a page is folded over so as to reach the left edge in such a way that one end point of the crease is on the right edge of the page and the other end point is on the bottom edge of the page as in the figure. If the width of the page is c inches, find the minimum length of the crease.
12. For $f(x) = 44 + 4x - 13x^2 + 18x^3 - 9x^4$ we give certain of the functional values and a graph on which the points of the table are shown. Sketch a graph through the points. Where do you think the maximum value of f is?



x	-2	-1	0	1	2	3
$f(x)$	-304	0	44	44	0	-304



13. Approximate the maximum value of the function

$$f(x) = 39 - 640x^2 - 1280x^3 - 640x^4$$

1-2. Area. The Integral.

The general concept of plane area is another of those geometrical ideas--like that of the direction of a curve at a given point--which remains elusive unless conceived in terms of limits. We already know a great deal about areas

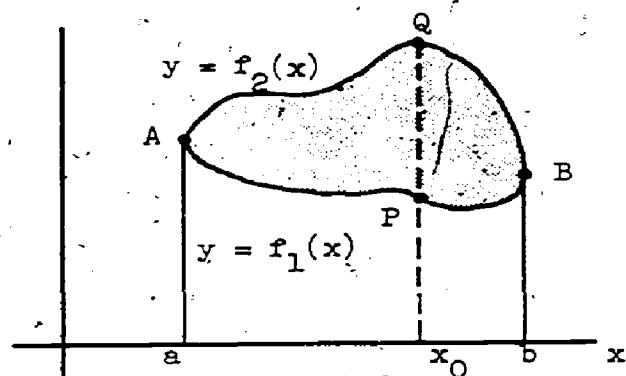


Figure 1-2a

from geometry. We know how to calculate the areas of triangles and, hence, the areas of all figures built up of triangles, that is, of all polygons. It is a problem of another kind to determine the area of a region with curved boundaries like the shaded region in Figure 1-2a.

You may know about the ancient Greek approach to the problem of determining the area of a circle: the area is described as the limit of areas of regular inscribed and circumscribed polygons. The Greeks were also able to calculate the areas of sections of a parabola--that is, regions bounded by a line segment and a parabolic arc. This was, in substance, the Greek contribution to the theory of area.

In the modern theory of area we successfully make general use of the basic idea of determining the area of a given region as the limit of approximations by polygonal regions. You may wonder, then, why the Greeks were not able to generalize this method. Historians usually attribute the limitations of the Greeks in this field to a lack of adequate general schemes for operating with numbers. It seems that the Greeks customarily thought of real numbers as magnitudes of geometrical objects rather than as entities studied independently of geometry. Nowadays we have a broad base for thinking both geometrically and numerically, and we take whichever point of view is the more convenient for the problem at hand. The enormous flexibility of this dual approach will enable you to solve handily problems which would have baffled the greatest Greek mathematicians.

To turn the geometrical description of the problem into a numerical one we introduce a coordinate system in the plane. For simplicity we place the axes so that the region in question is contained in the upper half-plane, $y \geq 0$, as in Figure 1-2a. Next we attack the problem of describing the region numerically. We know that some curves are the graphs of functions, and we are thus led to think of describing the boundary curve in terms of functions. The only difficulty is that the boundary curve is closed, so that a vertical line will generally meet the curve more than once. In Figure 1-2a, the vertical

Line $x = x_0$ meets the curve in two points P and Q. For this special curve the boundary curve can be divided into a lower arc APB and an upper arc AQB so that a vertical line intersects each arc no more than once. Each arc can then be considered as the graph of a nonnegative function defined on the interval $[a, b]$, that is, defined for all values of x satisfying $a \leq x \leq b$, where a is the abscissa of A and b the abscissa of B. The numerical description of the boundary curve is now given in terms of two functions, a lower function $f_1 : x \rightarrow f_1(x)$ corresponding to the arc APB and an upper function $f_2 : x \rightarrow f_2(x)$ corresponding to AQB. Since we have two functions to deal with, we are led to separate the calculation of the area into two parts. The area we seek is simply the difference between the areas of two regions of the same type. These are regions cut out of the strip between the vertical lines through A and B: the smaller region is bounded above by the graph of f_1 and below by the x-axis; the larger region is bounded above by the graph of f_2 and below by the x-axis (Figure 1-2a).

We have reduced the problem of determining the area of the given region to the problem of determining the area of regions of a certain standard type, regions describable in terms of a single function. Of course, the region we began with was especially simple. In more complicated cases a vertical line may meet the boundary curve in more than two points and we shall then need more than two functions to describe the curve (Figure 1-2b). We can still approach the problem by introducing standard regions, one for each function; the method for doing so in general is left for you to think about since the details are not relevant at the moment.

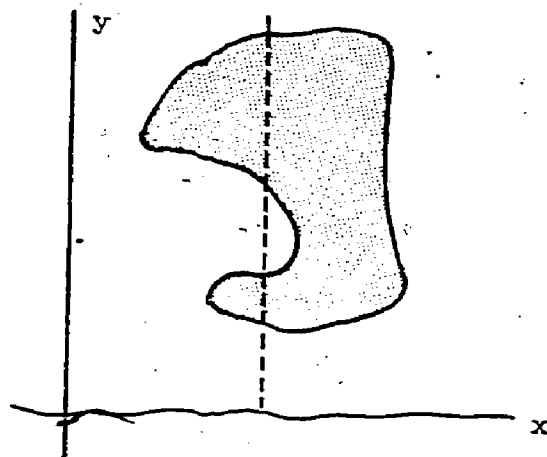


Figure 1-2b

We are left with the problem of calculating the area of a standard region, for example, the shaded region in Figure 1-2c.

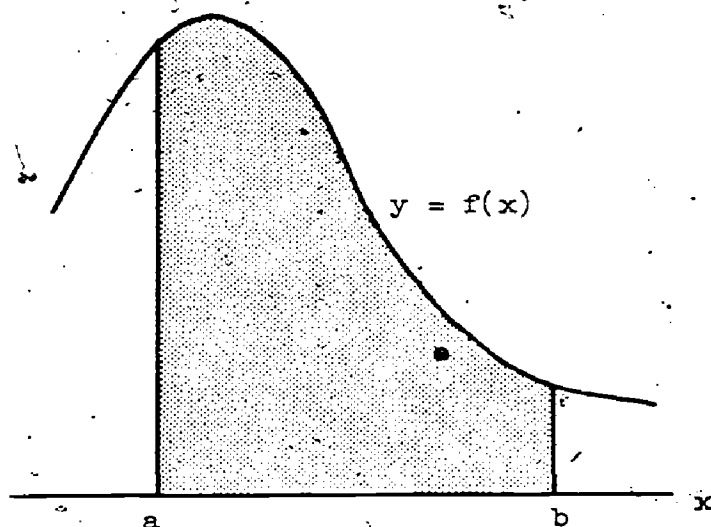


Figure 1-2c

In general, given a nonnegative function f on an interval $[a, b]$, we define the corresponding standard region as the set of points (x, y) for which $a \leq x \leq b$ and $0 \leq y \leq f(x)$, (Figure 1-2c). Again we are faced with the problem of determining a number, the area of the standard region based on the interval $[a, b]$, and the problem is apparently insoluble by any of the old methods unless the graph of f is composed of line segments. Again we

approach the problem by treating the area as a limit. We approximate the area by polygonal areas as the Greeks did, but we are looking for a systematic scheme of approximation, one that does not depend on the particular function involved.

A first crude estimate of the area A of the standard region on the interval $[a, b]$ can be given in terms of the minimum value m and the maximum value M of $f(x)$. (See Figure 1-2d.) Clearly, the rectangle of height m based on the interval $[a, b]$ is contained in the given standard region; the given region, in turn, is

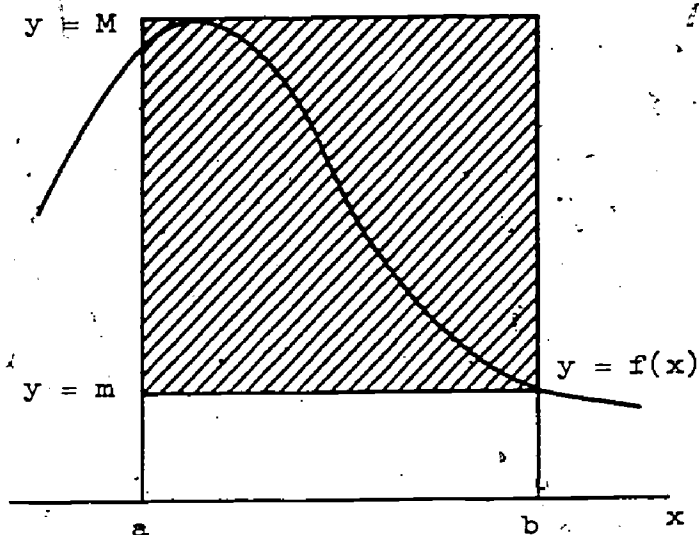


Figure 1-2d

contained in the rectangle of height M on the same base. For the area A of the region we then have estimates from below and above:

$$m(b - a) \leq A \leq M(b - a).$$

If we approximate A by either of these estimates or by any value in between, then we cannot be in error by more than $(M - m)(b - a)$; that is, by the area of the hatched region indicated in Figure 1-2d.

This simple method of estimation can easily be refined in a straightforward way. For this we only have to observe that the minimum m^* of $f(x)$ on any subinterval $[x_1, x_2]$ cannot be less than the overall minimum m (see Figure 1-2e); similarly the maximum M^* of $f(x)$ on the same subinterval cannot be greater than the overall maximum M ; that is,

$$m \leq m^* \text{ and } M^* \leq M.$$

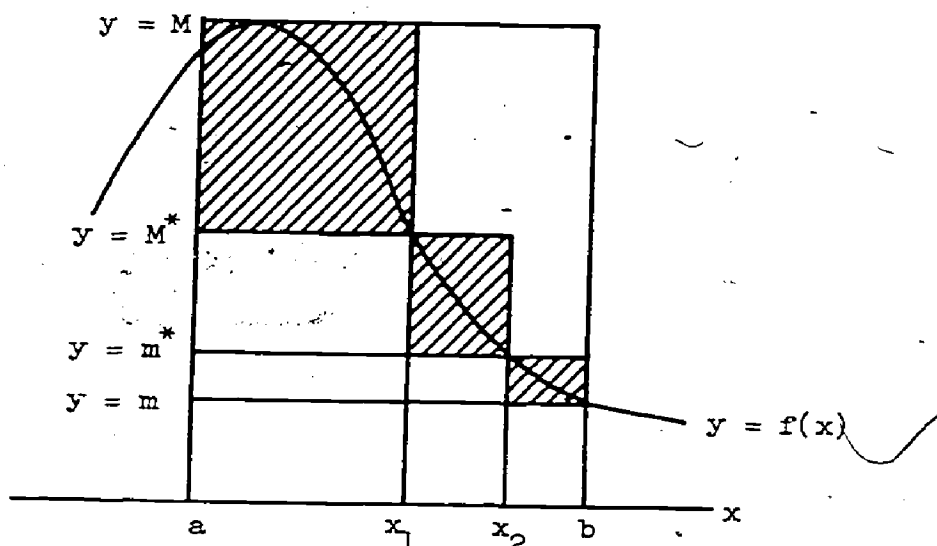


Figure 1-2e

It follows for the area A^* of the standard subregion based on the interval $[x_1, x_2]$ that

$$m(x_2 - x_1) \leq m^*(x_2 - x_1) \leq A^* \leq M^*(x_2 - x_1) \leq M(x_2 - x_1).$$

From this we see that the largest possible error in estimating the area of the subregion on $[x_1, x_2]$ has been reduced from the former value of $(M - m)(x_2 - x_1)$ to $(M^* - m^*)(x_2 - x_1)$.

This suggests that we can reduce the maximum error in estimating the entire area A by subdividing the interval $[a, b]$ into smaller intervals and making the same sort of estimate of area separately for each subregion. The sum of these estimates will be a better approximation to the area A than the

first, crude estimate. For the subdivision of Figure 1-2e this means reducing the margin of error from the area of the hatched region in Figure 1-2d to the area of the hatched region in Figure 1-2e.

The process of subdividing the interval $[a,b]$ can be repeated indefinitely, and this suggests that we now try to reduce the maximum error below any given tolerance by making the subdivision fine enough. If we do so the area will be given as the limit of both upper and lower estimates. The general limiting process by which we have characterized area, in particular, is called integration, and the corresponding limit is called an integral.

A good way to see how this general approach works, is to try it out on a specific function. For this purpose we try to find the area of the shaded region in Figure 1-2f. That is, we try to evaluate the integral A of $f : x \rightarrow \sqrt{x}$ from 0 to 1.

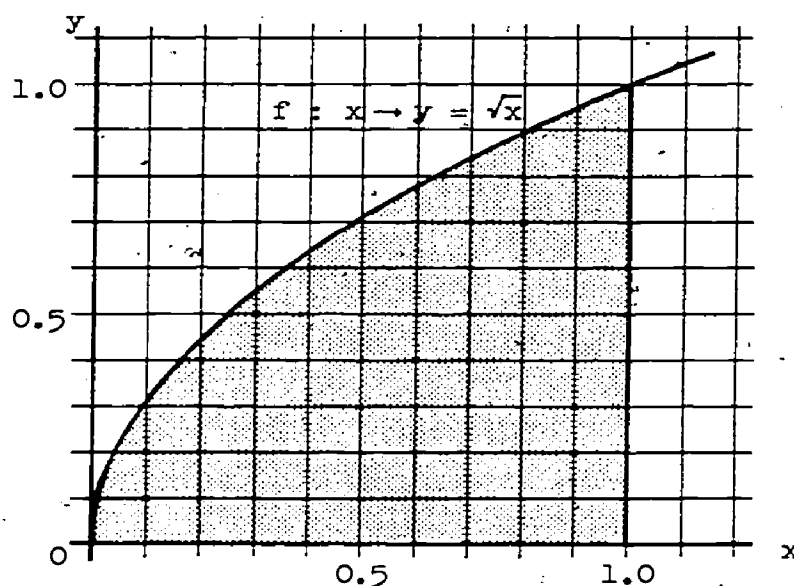
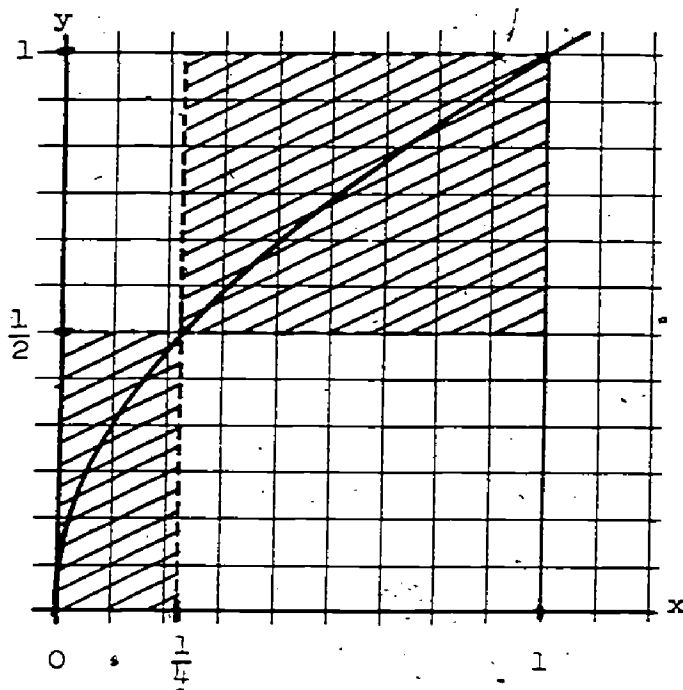


Figure 1-2f

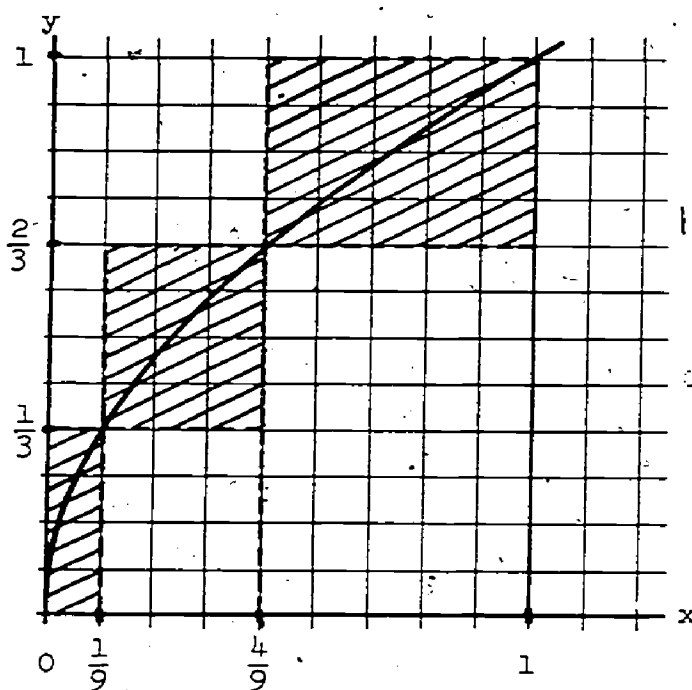
For this function we observe that $f(x)$ increases as the value of x increases. It follows that in any interval the minimum value of $f(x)$ occurs at the left endpoint and the maximum value at the right. Hence, for the entire interval $[0,1]$, the minimum value of the function is $f(0) = 0$ and the maximum is $f(1) = 1$. The areas of rectangles with these values as heights and a base of 1 are, respectively, 0 and 1. Thus, as a preliminary estimate, we know that the area A we seek is between 0 and 1.

In order to refine the estimate, we subdivide the base into n parts and denote the successive endpoints of the n subintervals by x_0, x_1, \dots, x_n where $0 = x_0 < x_1 < \dots < x_n = 1$. For simplicity in dealing with $f(x) = \sqrt{x}$,

we choose the endpoints so that $\sqrt{x_1} = \frac{1}{n}$, $\sqrt{x_2} = \frac{2}{n}$, $\sqrt{x_3} = \frac{3}{n}$; ...; that is, $x_1 = (\frac{1}{n})^2$, $x_2 = (\frac{2}{n})^2$, $x_3 = (\frac{3}{n})^2$, ... We illustrate this procedure for $n = 2$ and $n = 3$ in Figure 1-2g, and then calculate lower and upper estimates of the integral A for each of these cases by adding areas of rectangles. The area of the hatched region in each figure is the difference between the upper and lower estimates of the area A and represents the margin of error for the particular subdivisions.



$n = 2$



$n = 3$

Figure 1-2g

For $n = 2$,

$$0 \cdot \frac{1}{4} + \frac{1}{2} \cdot \frac{3}{4} \leq A \leq \frac{1}{2} \cdot \frac{1}{4} + 1 \cdot \frac{3}{4}$$

or

$$\frac{3}{8} \leq A \leq \frac{7}{8}$$

We have reduced the maximum error of 1 obtained without subdivision ($n = 1$) to $\frac{7}{8} - \frac{3}{8} = \frac{1}{2}$.

Taking $n = 3$, we obtain similarly

$$0 \cdot \frac{1}{9} + \frac{1}{3} \cdot \frac{3}{9} + \frac{2}{3} \cdot \frac{5}{9} \leq A \leq \frac{1}{3} \cdot \frac{1}{9} + \frac{2}{3} \cdot \frac{3}{9} + 1 \cdot \frac{5}{9}$$

or

$$\frac{13}{27} \leq A \leq \frac{22}{27};$$

the maximum error is reduced to $\frac{9}{27} = \frac{1}{3}$.

In general, for n subdivisions we have

$$A \geq \frac{0}{n} \left[\left(\frac{1}{n} \right)^2 - \left(\frac{0}{n} \right)^2 \right] + \frac{1}{n} \left[\left(\frac{2}{n} \right)^2 - \left(\frac{1}{n} \right)^2 \right] + \frac{2}{n} \left[\left(\frac{3}{n} \right)^2 - \left(\frac{2}{n} \right)^2 \right] \\ + \dots + \frac{n-1}{n} \left[\left(\frac{n}{n} \right)^2 - \left(\frac{n-1}{n} \right)^2 \right]$$

or

$$A \geq \frac{1}{n^3} \left[0 \cdot 1 + 1 \cdot 3 + 2 \cdot 5 + \dots + (n-1)(2n-1) \right]$$

Similarly,

$$A \leq \frac{1}{n} \left[\left(\frac{1}{n} \right)^2 - \left(\frac{0}{n} \right)^2 \right] + \frac{2}{n} \left[\left(\frac{2}{n} \right)^2 - \left(\frac{1}{n} \right)^2 \right] + \frac{3}{n} \left[\left(\frac{3}{n} \right)^2 - \left(\frac{2}{n} \right)^2 \right] \\ + \dots + \frac{n}{n} \left[\left(\frac{n}{n} \right)^2 - \left(\frac{n-1}{n} \right)^2 \right]$$

or

$$A \leq \frac{1}{n^3} \left[1 \cdot 1 + 2 \cdot 3 + 3 \cdot 5 + \dots + n(2n-1) \right]$$

Taking the difference between these results, we find for the maximum error E_n for this subdivision that

$$E_n = \frac{1}{n^3} \left[1 \cdot 1 + 1 \cdot 3 + 1 \cdot 5 + \dots + 1 \cdot (2n-1) \right]$$

or

$$E_n = \frac{1}{n^3} \left[1 + 3 + 5 + \dots + (2n-1) \right]$$

The expression in braces is an arithmetic progression for which we know the sum. We obtain, finally

$$E_n = \frac{1}{n^3} \cdot n^2 = \frac{1}{n}.$$

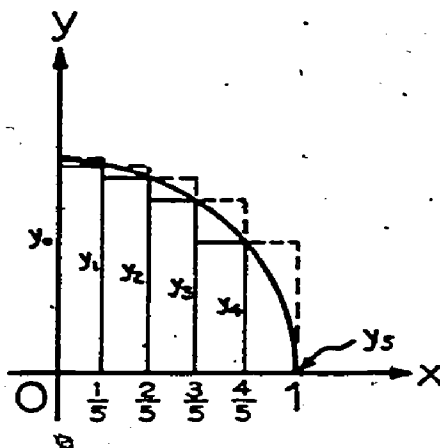
For this method of subdivision, then, we can reduce the error below any given tolerance simply by taking n big enough: given an error tolerance a , we simply take $n > \frac{1}{a}$.

It may seem that we have not solved the problem of finding the number A . All we know is that we can approximate A to within any given margin of error. Nonetheless, in describing the integral A of f from 0 to 1 as the limit of a set of approximations we have left no room for ambiguity. We still may feel cheated. We would like to have a familiar representation for A like, say, $A = \frac{2}{3}$ (which, by the way, is exactly what A is in this case). Later we shall see how to obtain such a number in cases like this. Still, it is important to know that we cannot always expect the solutions of our problems to take a familiar form; often the simplest and most comprehensible description of a number is its description as the limit of a set of approximations.

Exercises 1-2

1. Sketch the graph of $f : x \rightarrow \frac{1}{x^2}$ from $x = 1$ to $x = 2$.
- (a) Estimate the integral A of f from 1 to 2 using any method that you wish.
- (b) If you use the method of Section 1-2 and approximate A by rectangular areas, what are upper and lower estimates for two subintervals? What is the maximum error E_2 for this subdivision?
- (c) Using the same procedure, subdivide the base into n equal parts and find the maximum error E_n for this subdivision.
- Hint: Denote the successive endpoints x_0, x_1, \dots, x_n where $x_0 = 1 = \frac{n}{n}$, $x_1 = 1 + \frac{1}{n} = \frac{n+1}{n}$, $x_2 = 1 + \frac{2}{n} = \frac{n+2}{n}$, \dots , $x_n = 1 + \frac{n}{n} = \frac{2n}{n}$.
- (d) How can the maximum error E_n be brought below the error tolerance 0.03?
- (e) How can the maximum error E_n be brought below any given error tolerance?

2. A circle of unit radius has an area of π square units. Consequently, we can approximate π by using the same method as in No. 1. Here we shall use a quarter of a circle and obtain an approximation for $\frac{\pi}{4}$. If we use five equal subdivisions the rectangles enclosed in the desired region will have heights of y_1, y_2, y_3, y_4 , and y_5 while the rectangles containing the desired region will have altitudes of y_0, y_1, y_2, y_3, y_4 , where y_i is the ordinate of the quarter-circle $y = \sqrt{1 - x^2}$ at $x = \frac{i}{5}$, $i = 0, 1, \dots, 5$. The base of each rectangle is $\frac{1}{5}$. So



$$\frac{1}{5} (y_1 + y_2 + y_3 + y_4 + y_5) < \frac{\pi}{4} < \frac{1}{5} (y_0 + y_1 + y_2 + y_3 + y_4).$$

The difference between these larger (ℓ) and smaller (s) estimates is

$$\frac{y_0 - y_5}{5} \text{ or } \frac{1}{5}$$

Then the average of the larger and smaller estimates can only differ from $\frac{\pi}{4}$ by less than half this difference; that is,

$$\left| \frac{\pi}{4} - \frac{\ell + s}{2} \right| < \frac{1}{10}$$

How many equal subdivisions should one take to obtain an approximation of π correct to within $\pm .01\%$ of π ?

1-3. The Scope of Calculus

Calculus is the study of the derivative and the integral, the relationship between these concepts and their applications. The derivative and the integral may be interpreted geometrically as slope and as area, but these are only two among a wide range of interpretations and applications. We emphasized slope and area in the preceding sections in order to introduce parts of the subject in an intuitive, geometrical way. However, the concepts of derivative and integral are universal, and their incorporation into a calculus, a system of reckoning, enables us to solve significant problems in all branches of science. The intuitive approach we followed to this point is useful and suggestive, but it needs to be and will be supplemented by more careful study that indicates the broad range of application of these methods as well as their limitations. Before starting our systematic development of the subject, we want to emphasize the universality of the concepts of derivative and integral. We also want to stress that the problems considered in the initial chapters are primarily vehicles for the development of theory; they do not begin to suggest the full scope of applications of the calculus.

An intuitive geometrical introduction must necessarily be based on very familiar steps. However, the steps are so familiar that it may not seem that they could lead to entirely new methods for solving completely different problems than those encountered in earlier courses--problems that range from the spreading of rumors to the motions of the planets. No methods of the earlier courses would help you to answer questions such as: How did you first hear about calculus? How did you first hear about Helen of Troy? To frame such questions mathematically, we must first isolate some essential features: Some stories spread like diseases; others die out. If the story is too dull, nobody bothers to repeat it. But if the story is good, some of the people who hear it (and remember it) pass it along. Starting from these ideas, how far could you get by familiar methods?

The same concepts that are basic to the spreading of stories are also basic to the processes of forgetting and learning. So many of the facts we stuffed into our heads to pass a test and never used afterwards seem to have vanished. Others, that we met repeatedly and actually worked with have become so much a part of ourselves that we can hardly remember not having known them. Our first exposures to these facts may not have taken, but repeated encounters in different contexts finally left their mark.

Can we find a scientific explanation for the way that stories spread--the way that we learn, and forget, facts? Starting from appropriate assumptions

(the mathematical model) we can and later shall describe some aspects of such processes with the aid of the calculus. Such processes illustrate a broad class of phenomena whose unifying features are the basic mathematical models for growth, decay, and competition. Besides helping to describe the spreading of rumors, and learning or forgetting, these same mathematical models serve to clarify such natural phenomena as radioactive decay, the attenuation of sunlight by a cloudy sky, the progress of chemical reactions, the growth of bacterial colonies, or the spread of disease through a city. In each of these situations, the essential feature is that the amount of some quantity is changing (with respect to time, or distance, or whatever) at a rate proportional to the amount already present. A process of this sort can be mathematically described by a certain type of equation (a differential equation) whose solutions, at least in the simplest cases, are combinations of exponential functions.

Other processes of nature change in a cyclic or periodic way; they repeat in identical form each year, each day, perhaps each second. The planetary motions, the tides, the daily routine of life; the harmonious chords of music, the propagation of X-rays through crystals--even the colored sheen of butterfly wings--all depend on periodic phenomena. For such processes, the rate of change of the rate of change of some quantity is (negatively) proportional to the quantity itself, and our mathematical model leads to a different class of differential equations whose solutions, in the simplest cases, are combinations of trigonometric functions.

With the calculus we may also investigate more complicated natural processes that involve a combination of growth or decay with some sort of cyclic behavior. We may also solve much simpler problems: How much time will it take to drive 1000 miles if you start at a speed of 30 miles per hour, but increase your speed by 3 miles per hour for each hour of driving? At what angle should you throw a ball for it to travel as far as possible? In what directions with respect to the sun are the rainbow colors strongest?

These problems and many others which the calculus solves involve rates of change; this is the province of the differential calculus. A second broad variety of questions is concerned with totality--the summing of small effects; this is the province of the integral calculus. By recording your speed during a long trip, with many accelerations and decelerations, can you calculate how far the trip has taken you? If we know how a single drop of ink spreads on a blotter, can we predict what happens if we spill the whole bottle? With the differential calculus we can construct a mathematical model for radioactive processes; how do we go about shielding a nuclear reactor? To explain the

sunburn we acquired on a cloudy day, how much light reaches us from the entire overcast sky and by reflection from our surroundings?

Such summation or integration problems are closely related to the rate-of-change or differentiation problems: the total effects result from an addition of small variations. Therefore we do not study separately an integral calculus and a differential calculus. We study a calculus comprising both differentiation and integration, and each aspect helps us to understand and apply the other.

In this course we strongly emphasize applications, not only to show that calculus provides useful methods and concepts for the sciences, but because so much of the calculus was developed to solve specific problems. Most of the applications of mathematics we now consider differ from those of previous courses in that they emphasize the effects of variation or summation. "Calculus" was tailor-made to treat such problems. Except for the simplest problems of this type, the methods of arithmetic, geometry, and algebra are inadequate, and even for the simpler problems the methods of calculus are the more efficient.

The calculus was invented to treat problems of physics. As the calculus grew into the larger branch of mathematics known as analysis its range of application expanded enormously. To analysis we owe much of the progress in the physical sciences and modern engineering, and more recently in the biological and social sciences. The concepts and operations introduced by the calculus provide the right language and the right tools for the major part of the applications of mathematics to the sciences.

The great advance which takes the calculus beyond algebra and geometry is based on the concept of limit. Our initial examples were chosen, not merely because they were simple, but because they illustrate the essentials. The basic limit procedure of the differential calculus is typified by the problem of finding the slope of a curve; the basic limit procedure of the integral calculus is typified by the problem of finding the area enclosed by a curve. The slope is found as a derivative, the area as an integral, and superficially these appear to be unrelated. But there is only one calculus: derivative and integral are complementary ideas. If we take the slope of the graph of the area function, we are brought back to the curve itself. If we take the area under the graph of the slope function, we find the original curve again. The limit concept, in its guises of derivative and integral, together with this inverse relation between the two, provides the fundamental framework for the calculus.

In this first chapter, we have introduced the calculus by elementary examples that motivate the study of derivative and integral. Having started intuitively, we go on to indicate the mathematical problems that arise when we try to give a precise foundation to our intuition; we solve these mathematical problems, develop methods for obtaining solutions to basic problems involving variation and summation, and then apply these solutions systematically to illustrate the interrelations of the calculus and the sciences.

In our approach, we have tried to maintain a balance of topics and of viewpoints that will meet the requirements of students who will become mathematicians, others who will become scientists primarily interested in applications, and still others for whom mathematics will become simply one of many deep intellectual experiences during their education. For the student of science, a fluent intuitive grasp of the subject may often seem to be his primary need; for the student of mathematics, it may seem that a careful deductive development is more essential. These different viewpoints necessarily conflict on occasion, but more frequently they supplement each other in providing an overall picture of this new mathematical subject. Both the scientist and the mathematician gain by a complete command of both points of view, and we regard this as an ideal worth striving for; the serious student should carry away a deeper appreciation of both these views and of their interrelations.

Historically, the replacement of an intuitive basis for the calculus (the method of infinitesimals) by a careful logical structure (the method of limits) marked a vital phase in the development of mathematics. This phase is far from complete. We are still learning how to combine the inspired use of intuition in approaching new problems with the effective use of logic, not only to verify our intuition, but to permit generalizations of broader applicability. Today, most mathematicians appreciate the essential roles of both intuitive and deductive procedures, not only for creating the calculus but for learning it, and we have tried to make both equally available to you.

In summary, we wish to show you how the effort to solve important problems leads to methods of the calculus; how the attempt to make the best use of these methods and to understand their full scope leads to the development of the calculus as an independent study; and how the products of this study in turn lead to deeper insight into the original problems. Just as science enriches mathematics by providing concrete models and significant problems, mathematics enriches science by providing system and organization, as well as solutions of problems.

Exercises 1-3

The following is a quotation from G. Polya: "If you cannot solve the proposed problem, try to solve first some related problem."* In the problem of Section 1-1 concerning the shipping of books, a "related problem" is considered when we switch our unfocused interest in the problem of finding the dimensions of a box with maximum capacity to the specific related problem of determining properties of the function f given by

$$f(x) = V = x^2(72 - 4x) .$$

A final related problem is considered when looking for the point of the graph of this function at which the slope is zero.

What are the related problems in

- (a) Ex. 1-1, number 5?
- (b) Ex. 1-1, number 7?
- (c) Ex. 1-1, number 8?
- (d) Ex. 1-1, number 10?
- (e) Ex. 1-2, number 1?

*Polya, G., How To Solve It, 2nd ed., New York: Doubleday, 1956; p. 10.

Chapter 2

THE IDEA OF DERIVATIVE

2-1. Introduction.

Although intuitive geometrical thinking may serve as a guide to the creation of a mathematical concept or proof, it is not always reliable. There is the danger of reading properties of a picture into a precisely formulated mathematical concept when those properties cannot be logically derived.

Example. Consider the function f , defined by

$$f(x) = \begin{cases} 1, & x > 1 \\ x, & -1 \leq x \leq 1 \\ -1, & x < -1 \end{cases}$$

The graph of f is exhibited in Figure 2-1 in three different scales where the unit is

- (a) 1 inch,
- (b) $\frac{1}{8}$ inch,
- (c) 10^{-6} inch.

The slope of the graph of f at $x = 0$ must certainly be defined as the slope of the line given by $y = x$; thus the slope at $x = 0$ is 1. Yet in Figure 2-1c the graph of f is not perceptibly distinguishable from the x -axis, which has the slope 0.

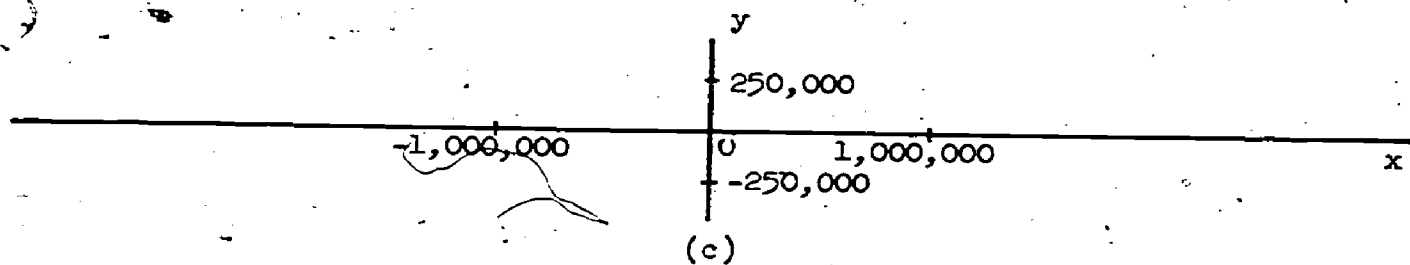
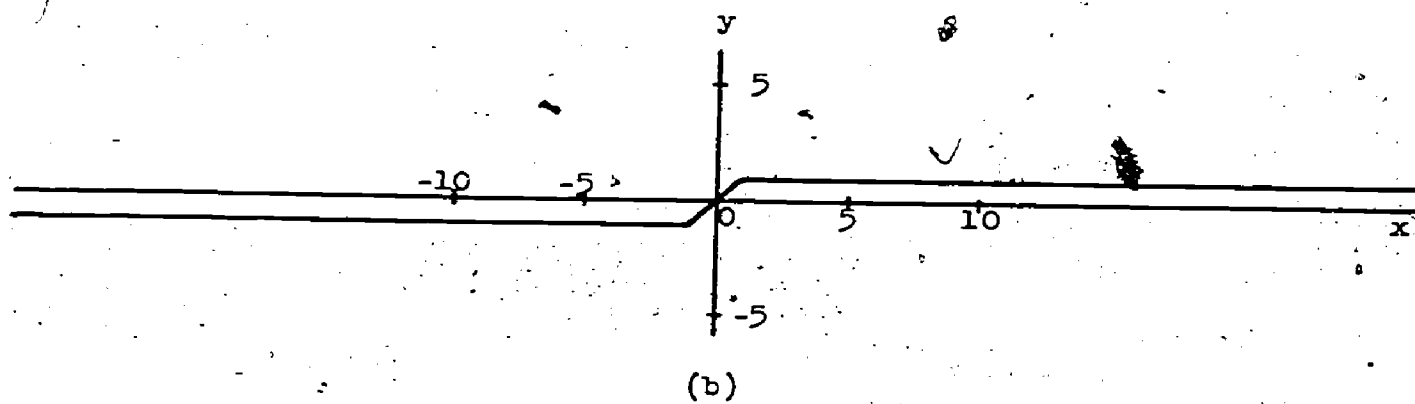
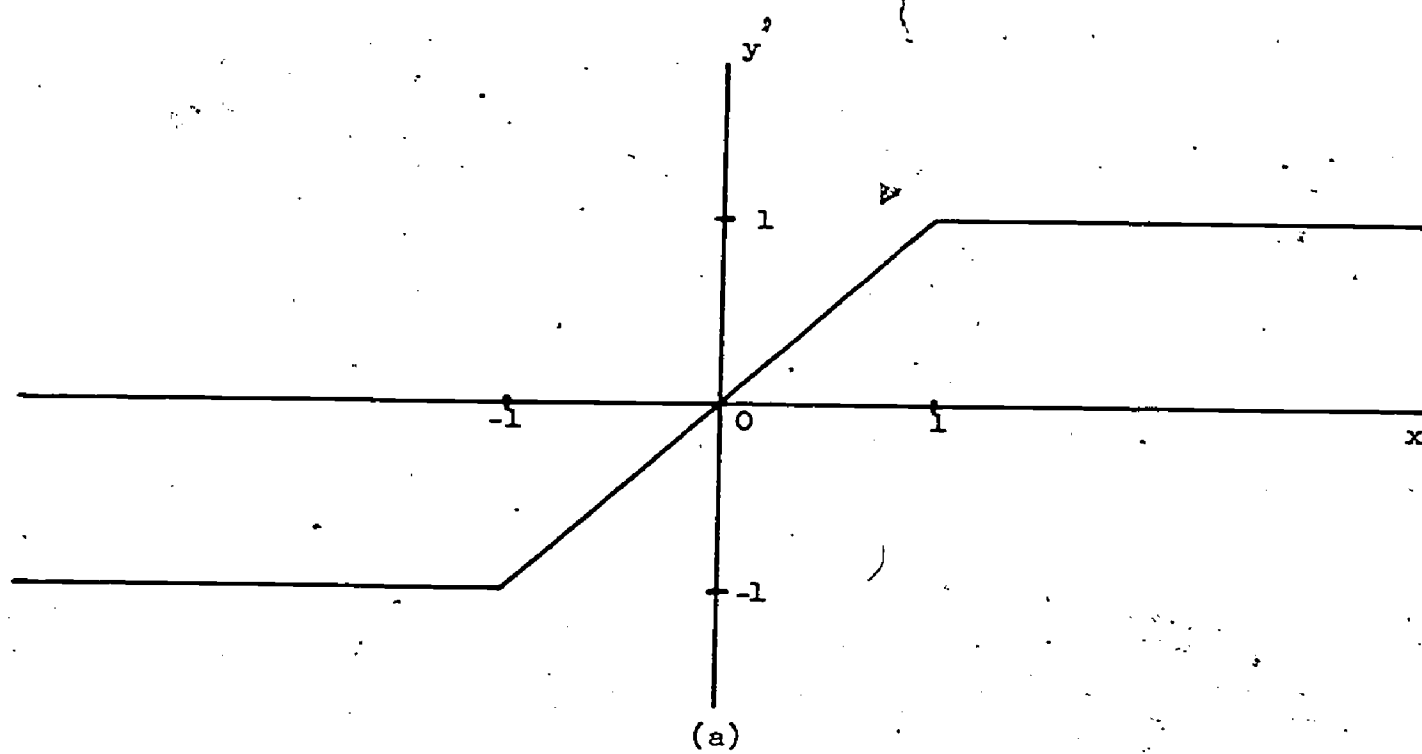


Figure 2-1

If, as in this example, an accident such as a particularly unfortunate choice of scale can obscure exactly those features which we are trying to explore, it makes sense to try to put our ideas in a form which enables us to reach conviction by precise logical argument rather than by appeal to our fallible perceptions. For this reason, without rejecting geometrical intuition as a guide in directing our thoughts, we shall express our concepts analytically in terms of number and take numerical concepts as the basis of proof in all that follows.

The treatment of the parcel post problem in Section 1-1 shows in a particular case how the geometrical concept of direction for a curve can be given a precise analytical interpretation in terms of slope. The slope of a curve at a given point is a local property in the sense that it is completely determined by any arc containing the point, no matter how small. For the purpose of describing such a local property, the idea of limit is especially appropriate. The slope of a curve at a point is defined as the limit of a set of slopes of chords to the curve. In this chapter we shall introduce this limit in terms of a purely analytical concept, the derivative of a function at a point. The analytical concept, divorced from its geometrical interpretation as slope, will be seen to have other realizations. Next, for some simple examples we shall see how to compute such limits, but we shall not treat the general computational problem until we have developed a clearer understanding of limit (Chapter 3).

It should not be thought that the methods of geometry are entirely powerless in the present context. The ancient Greeks treated the problem of defining the direction of a curve at a point by finding the tangent line at the point, the tangent being the line through the point which has the direction of the curve there. They were able to construct the tangents to all the conic sections (circle, parabola, ellipse, hyperbola) and even to such complicated curves as the spiral of Archimedes. In the end, though, the Greeks were unable to solve the problem of drawing tangents to more than a limited class of curves whose special geometrical properties made the problem tractable. What limited them was the fact that they had no general way of defining a curve, say, in terms of functions; that had to wait until the invention of analytic geometry by Descartes.

2-2. Slope.

The idea of slope for a curve at one of its points is a generalization of the idea of slope for a straight line. In coordinate-geometry, the direction angle or angle of inclination of a straight line is defined as the angle measured in the counterclockwise sense from any horizontal line given by $y = \text{constant}$. (Figure 2-2a.) Although direction angle is the more intuitive geometric concept, it is simpler analytically to work with the concept of slope, the tangent of the direction angle. For a straight line the slope m is defined by any two of its distinct points (x_1, y_1) and (x_2, y_2) , namely

$$m = \frac{y_2 - y_1}{x_2 - x_1},$$

where the value m given by this formula is independent of the choice of the two points and also of the order in which they are given. We remind you that the slope defines the direction angle unambiguously, and therefore serves adequately as a numerical characterization of the angle.

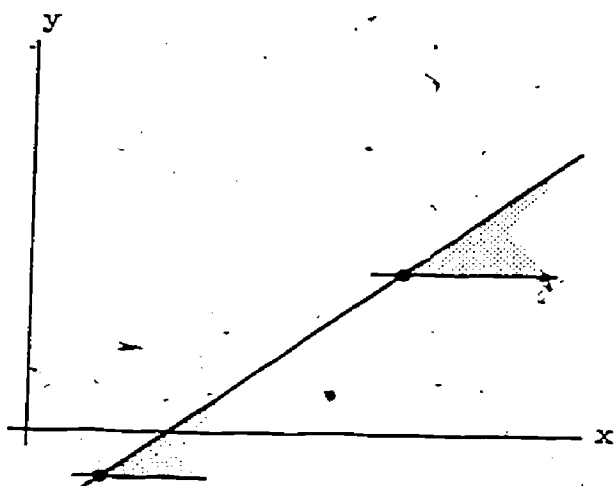


Figure 2-2a

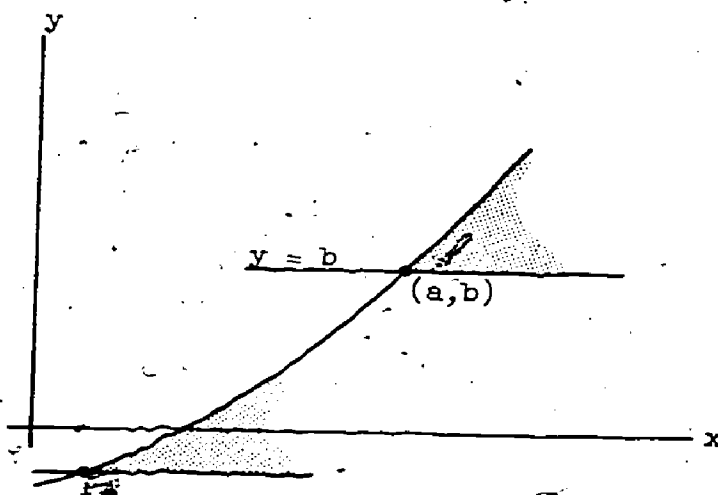


Figure 2-2b

For a curve that is not a straight line the direction angle changes from point to point (Figure 2-2b), and we must define this angle, and therefore the slope, separately at each point. We suppose that an arc of the curve containing the point is the graph of a function f . Following the idea of the example in Section 1-1, we expect the slope of the graph at the point (a, b) , where $b = f(a)$, to be approximated by the slopes of chords joining (a, b) to nearby points (x, y) on the graph, where $y = f(x)$ (Figure 2-2c). If these chordal slopes are adequate approximations to a single number m , in the sense that it is possible to make the error in the approximation smaller

than any prescribed tolerance by a choice of x sufficiently near to a , then we accept m as the slope of the curve at the point (a,b) . In other words, we define the slope at (a,b) as the limit of the set of approximations furnished by the slopes of chords to nearby points.

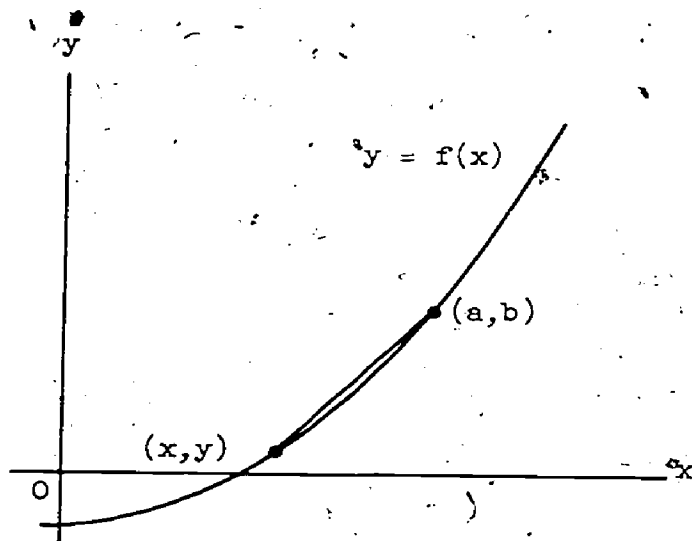


Figure 2-2c

To see how this idea of approximation works in practice, let us try it for the simplest nonlinear graph we know, the curve with equation $y = x^2$, at the point $(1,1)$. For the chord joining the points $(1,1)$ and (x,y) , we calculate the slope

$$r(x) = \frac{y - 1}{x - 1}$$

for a number of values of x taken successively closer to 1, and obtain the table

x	1.1	0.99	1.001	0.9999	1.000001
$r(x)$	2.1	1.99	2.001	1.9999	2.000001

From this table it is hard to escape the conclusion that the values of $r(x)$ are successively closer to 2 and that the error of approximation remains within any prescribed tolerance if x is close enough to 1. We may well believe that the slope m of the curve with equation $y = x^2$, at the point $(1,1)$, is $m = 2$.

Belief on the basis of a critical appraisal of the evidence is an excellent way to find the correct approach to a mathematical question; but no mathematical investigation is complete until the belief is substantiated by deduc-

tive, logical proof. To prove that the slope is actually 2 we must show that the error in approximating 2 by the slope $r(x)$ of the chord joining $(1,1)$ to (x,y) can be controlled by a choice of x sufficiently close to 1. In numerical terms, the error in approximating 2 by the slope $r(x)$ is $|r(x) - 2|$, and for any prescribed tolerance $\epsilon > 0$ we must be able to ensure that

$$(1) \quad |r(x) - 2| < \epsilon$$

if the distance between x and 1 is small enough. We have

$$r(x) = \frac{y - 1}{x - 1} = \frac{x^2 - 1}{x - 1} = \frac{x - 1}{x - 1} (x + 1).$$

The distance between x and 1 is defined as the absolute value of their differences, namely $|x - 1| = |1 - x|$. See Appendix 1-3.

We stress the fact that $r(x)$ has no meaning when $x = 1$; when $x = 1$, the denominator of the expression for $r(x)$ is zero. However, since $r(x)$ is defined as the slope of a chord, the value $x = 1$ is excluded from the domain of r . In that domain the ratio $\frac{(x - 1)}{(x - 1)}$ which appears in the last expression for $r(x)$ is constant and has the value 1, so that

$$r(x) = x + 1 \quad (x \neq 1).$$

It follows that

$$|r(x) - 2| = |x - 1|.$$

We see at once that by taking the distance between x and 1 smaller than ϵ we can guarantee that the inequality (1) holds.

Because we wished to exhibit clearly the procedure for verifying that a given value m is the slope of the graph of f at one of its points, we have kept the preceding example very simple. In certain respects it is too simple. In the first place, we assumed that the curve had a slope m , i.e., that the slopes of the chords had a definite limit. Then we were able to satisfy ourselves that such a limit exists by actually exhibiting its value $m = 2$ and showing that it is, in fact, the limit of the slopes of chords. In general, this procedure poses a second difficulty; we must find the numerical value of the slope before we can prove that it is the appropriate limit. How

does one find this value for a less simple function? Certainly the idea of discovering the value by inspecting a table of approximations is not always practical, as we may see from Exercise 2 below.

Finally, the choice of the point $(1,1)$ is too special for our purposes. It will be better to investigate the slope of the curve $y = x^2$ at any point (a,b) where $b = a^2$, without specification of a . For the slope $r(x)$ of the chord joining (a,b) to another point (x,y) of the graph we have

$$r(x) = \frac{y - b}{x - a} = \frac{x^2 - a^2}{x - a} = \frac{x - a}{x - a} (x + a),$$

where, of course, $x \neq a$. As before, the chordal slope $r(x)$ is defined for all real x except $x = a$, and so we have, on the domain of r , $\frac{x - a}{x - a} = 1$ and

$$(2) \quad r(x) = x + a.$$

We want to find the limit of $r(x)$ for values of x close to a , and it seems perfectly clear that this limit is $a + a = 2a$. In fact, if we accept $2a$ tentatively as the slope m of the graph at (a,b) , and set a tolerance of error ϵ , we can easily enforce the condition $|r(x) - 2a| < \epsilon$, for we have

$$|r(x) - 2a| = |(x + a) - 2a| = |x - a|.$$

We see that the error of approximation can be kept within the tolerance ϵ if we choose x so that $0 < |x - a| < \epsilon$. In words, we choose x at a distance less than ϵ from a , but we must exclude the choice $x = a$ because it defines no member of the set of approximations--the values of $r(x)$. We are sure now, not only that the curve $y = x^2$ has a definite slope m for $x = a$, but that we have also found its value, namely $m = 2a$.

In reviewing the argument just completed, we observe that (2) defines r as a linear function except at the one value $x = a$. The graph of r is thus a straight line with the one point $(a, 2a)$ deleted. It is intuitively reasonable to take $2a$ as the slope of $y = x^2$ at $x = a$. However, we reject intuition as our sole criterion, especially after considering examples like that of Figure 2-1. We insist on supporting our intuition with firm mathematical argument: we have verified that the value of the slope obtained on intuitive grounds agrees with the mathematical limit determined analytically.

In addition to the need for logical conviction, there is a more compelling reason for pursuing our argument in such detail. We wish to employ methods that are general, applicable to any function f and not merely to the very special case we have examined. In general, we may attempt to evaluate the slope of the graph of f at (a,b) , where $b = f(a)$, as the limit of the slopes

$$(3) \quad r(x) = \frac{y - b}{x - a} = \frac{f(x) - f(a)}{x - a}, \quad (x \neq a),$$

of chords joining (a,b) to nearby points (x,y) of the graph. (Again we observe that $r(x)$ is not defined at $x = a$.) Without knowing the function f explicitly, we cannot tell whether it is possible to proceed as before to obtain a simple expression for $r(x)$, valid on its domain, which will enable us to evaluate the slope. In fact, we shall give an example in Section 2-5 where such a simplification is impossible. For this reason, a treatment of slope as a limit must use the kind of argument that we have given.

Exercises 2-2

1. Find the slope of $y = r(x) = x^2 - 1$ at the point $(3, 8)$ by constructing a table of values of $r(x)$ for x successively closer to 3. Verify that your answer is the limit of the slopes of the chords.
2. For the function given by $f(x) = x^2 - x + 1$, tabulate the slopes of the chords joining $(\frac{5}{7}, f(\frac{5}{7}))$ to $(x, f(x))$ for $x = \frac{5}{7} + \frac{1}{10}, \frac{5}{7} - \frac{1}{10}, \frac{5}{7} + \frac{1}{100}, \frac{5}{7} - \frac{1}{100}, \frac{5}{7} + \frac{1}{1000}, \frac{5}{7} - \frac{1}{1000}$, etc., as far as your time, energies, and inclinations permit. Can you predict the limit of the approximations from an inspection of the table?
3. Each of the following curves passes through the origin.

(a) $y = f(x) = x^{1/3}$

(b) $y = f(x) = x^{2/3}$

(c) $y = f(x) = |x|^{3/2}$

For each curve construct a table listing the slopes of chords with one endpoint at the origin and the other at nearby points, $(x, f(x))$. As these points are taken successively closer to the origin (e.g., $|x| < 0.1$, $|x| < 0.01$), what information do you obtain about the slope of the curve? In your opinion, is it possible to define the slope of the graph at the origin? If so, what is the slope? If not, justify your answer.

4. (a) At each of the points $(1, 7)$ and $(2, 16)$, find the slope of $y = g(x) = 3x^2 + 4$ by constructing a table of values; then verify that your answer is the limit of the slopes of chords.
- (b) Find the slope of $y = g(x)$ at the point (a, b) on its graph, where $b = g(a)$.
- (c) Find the lowest point on the graph of g by methods of coordinate geometry.
- (d) Check your answer to (c) by using the result of (b). (If a hint is needed, one may be found in Section 1-1.)
5. (a) Find the slope of $y = x^3$ at (a, b) where $b = a^3$. (If a hint for the simplification of $r(x)$ is needed, one may be found in Section 1-1.)
- (b) Is there any lowest point on the graph of $y = x^3$? Is there any highest point? Is there any point where the graph is horizontal?

6. What is the relationship between the slopes of the function in Ex. 5 corresponding to the points $x = a$ and $x = -a$? Interpret this result graphically. Give examples of other functions having this property.
7. What is the relationship between the slopes of the function in Ex. 4 corresponding to the points, $x = a$ and $x = -a$? Interpret this result graphically. Give examples of other functions having this property.
8. (a) Find the slope of the graph of $h : x \mapsto 4x^3 - 3x^2$ at (a, b) , where $b = h(a)$.
- (b) Find all points where the graph of h is horizontal. Can you characterize these points as "highest" or "lowest," perhaps in a restricted sense?

2-3. General Quadratic Function.

We now apply the procedure of the last section to the class of all quadratic functions. Consider

$$f : x \longrightarrow Ax^2 + Bx + C,$$

where A, B, C are constants and $A \neq 0$. When $x = a$, the graph of f passes through the point (a, b) with $b = Aa^2 + Ba + C$, and we propose to calculate its slope there. Again we shall prove our result to be correct by showing that it is truly the limit of the approximations that define it, namely the slopes

$$r(x) = \frac{f(x) - f(a)}{x - a}$$

of chords through $(a, f(a))$ and $(x, f(x))$, with $x \neq a$. For the general quadratic function we obtain

$$\begin{aligned} r(x) &= \frac{Ax^2 + Bx + C - (Aa^2 + Ba + C)}{x - a} = \frac{A(x^2 - a^2) + B(x - a)}{x - a} \\ &= \frac{x - a}{x - a} [A(x + a) + B], \quad (x \neq a). \end{aligned}$$

When $x \neq a$, the function r coincides with the simple linear function $x \longrightarrow A(x + a) + B$. As before, we immediately guess that the slope at the point (a, b) must be $m = 2Aa + B$. We verify this guess by prescribing a tolerance ϵ and showing that the error of our approximations can be brought within the tolerance. For x sufficiently near to a the error in the approximation is $|r(x) - m|$, namely

$$|A(x + a) + B - (2Aa + B)| = |A(x - a)|.$$

In order to bring this error within the tolerance ϵ , we see that we must make a choice of allowable values of x depending on the coefficient A of the quadratic function. To ensure that we stay within the tolerance, that $|r(x) - m| < \epsilon$, we must restrict x to values for which $|x - a| < \frac{\epsilon}{|A|}$, so that

$$|r(x) - m| = |A(x - a)| = |A||x - a| < |A| \frac{\epsilon}{|A|} \leq \epsilon.$$

We have proved that our calculation of m gave us the correct slope--the

limit of the chordal slopes.

Example 2-3. Consider the parabola given by $y = f(x) = \frac{1}{2}x^2 + x + \frac{5}{2}$, where $A = \frac{1}{2}$, $B = 1$, and $C = \frac{5}{2}$. The lowest point on this graph may be found algebraically from the standard form $(x + 1)^2 = 2(y - 2)$ of the equation of the curve; it is the vertex $(-1, 2)$. Alternatively, using the idea introduced in the best-value problem in Section 1-1, we may find the lowest point directly as a point of zero slope. The slope at (a, b) is $m = (\frac{1}{2}) 2a + 1 = a + 1$, and it is zero for $a = -1$, where $b = f(-1) = 2$. (Note, however, that we do not as yet have a way to prove by the calculus that this is the lowest point; but the standard form does show that y cannot go below 2 on the graph.)

Exercises 2-3

1. Find the slope for $x = a$ of the general linear function $f : x \rightarrow Ax + B$ (where A and B are any constants except that $A \neq 0$) and compare your result to that obtained from the standard slope-intercept form of the equation of a straight line in coordinate geometry.

2. For what values of k does the line $y = k$ intersect the parabola

$$y = Ax^2 + Bx + C \quad (A \neq 0)$$

in

- (a) no points?
 - (b) 1 point?
 - (c) 2 points?
 - (d) What is the lowest or highest point of the given parabola?
3. (a) Find the highest point on the graph of

$$g(x) = 5 - 6x - x^2$$

using Ex. 2.

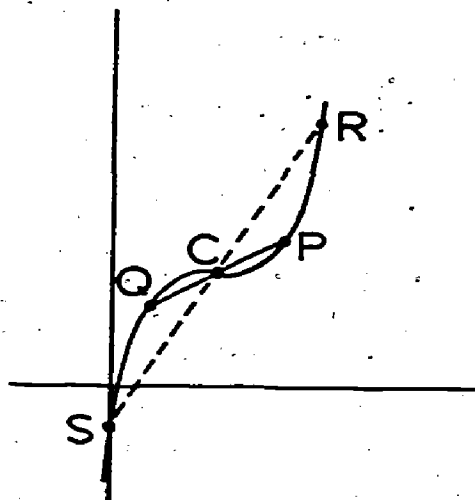
- (b) Explain geometrically why the point in (a) can also be obtained by finding where the slope of $g(x)$ is zero.
4. (a) What is the greatest possible number of points where the graph of a quadratic function $Ax^2 + Bx + C$ may be horizontal?
- (b) Is it possible for the graph to be horizontal at less than the maximum number of points, or nowhere horizontal? If the answer to either question is affirmative, give an example (in the form of a specific function).
5. (a) Given $y = f(x) = 20x - 3x^2$. Find the slope of the curve at the point (a, b) , where $b = f(a)$.
- (b) Where is the slope zero? How can you use this information in plotting the graph of f ?
6. (a) Find the slope of the curve with equation $y = h(x) = Ax^3 + Bx^2 + Cx + D$ (the graph of the general cubic function) at (a, b) , where $b = h(a)$; here A, B, C, D are any constants, except that $A \neq 0$. (If you need a hint for the simplification of $r(x)$, it may be found in Section 1-1.)
- (b) What is the greatest possible number of points where the graph of a cubic function h may be horizontal?

- (c) Is it possible for a cubic function to have its graph horizontal at less than ~~the~~ maximum number of points? If the answer is "Yes," give an example of such a function.
- (d) Is it possible for a cubic function to have its graph nowhere horizontal? If the answer is "Yes," give an example of such a function.
7. Show that the curve of Ex. 6 is centrosymmetric about the point $(-\frac{B}{3A}, h(-\frac{B}{3A}))$.

Centrosymmetric means that given the above point as a center, a line segment starting at any other point on the curve going through the center and extended will intersect the curve again so that the center is the midpoint of this segment.

$$PC = CQ \text{ and}$$

$$RC = CS.$$



8. A function h is defined by the formulas

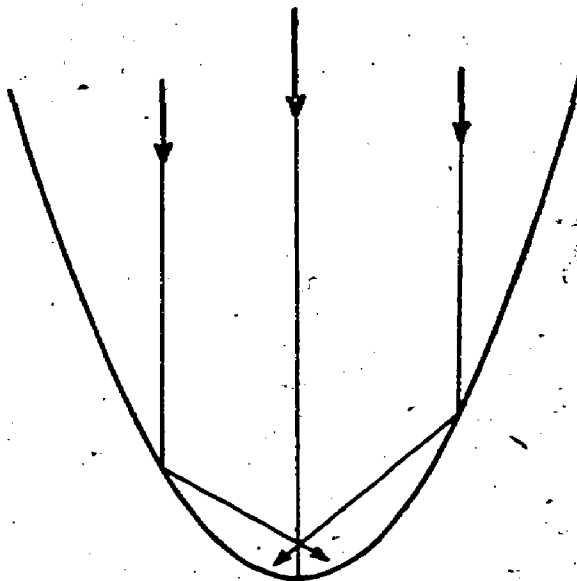
$$h(x) = \begin{cases} x^3 & , x < 1 \\ -x^2 + 2x & , x > 1 \end{cases}$$

Find the slope of the graph of h for $x < 1$ and for $x > 1$. Is it possible in your opinion to define a slope for the graph at $(1,1)$? Give an argument to support your answer. (A sketch may be helpful in answering the question.)

9. Until now our discussion of the idea of direction and slope for a curve has been generally at a theoretical level. Although we know from Section 1-1 that the concept of slope will ultimately be useful in "best-value" problems it is satisfying to have another more immediate application. You are probably familiar with the fact that large telescopes and automobile headlights use parabolic mirrors. A parabolic reflector can bring a bundle of parallel rays like those from a star* to a sharp focus. You are now able to demonstrate the sharp focusing property of the parabola.

* At stellar distances the deviation from parallelism of all rays reaching the earth from a given star is utterly negligible.

According to Heron's Law of reflection for a ray of light incident upon a smooth mirror, the incident ray and the reflected ray make equal angles with the mirror. Suppose the shape of the cross section through the axis of the mirror is given by the graph of $y = x^2$. Prove for all incident rays parallel to the y-axis that the reflected rays have a common point of intersection as in the figure. This common point is called the focus of the parabola. It can also be shown that this property characterizes the parabola; i.e., if the curve is such that all parallel rays pass through a common point after reflection, then the curve must be a parabola.



2-4. Velocity.

The kind of limit which is used to define the concept of slope is appropriate for other purposes. Now we shall see how it may be used to describe how rapidly an object is moving at a particular instant.

As a means of getting at this concept we compare the distance moved during a given time interval to the time elapsed; we introduce the average velocity: if the object is at a distance s_1 along its path at time t_1 and a distance s_2 along its path at time t_2 , then the average velocity is the rate or ratio $(s_2 - s_1)/(t_2 - t_1)$. Unfortunately, the average velocity tells us very little about the comparison of distance to time over any specific smaller time interval of the motion, or at any particular instant: knowledge that a motor trip was completed at an average velocity of 50 miles per hour gives no indication of the maximum speed or whether there were any stops.

We describe the course of the motion by measuring the distance s along the path from some reference point $s = 0$, setting a clock to the time $t = 0$ at some convenient instant of the motion (say its start), and then expressing the location of the object at time t by a function given by $s = \phi(t)$. Since a difference of values of s may be either positive or negative, the distance between two locations s_1 and s_2 is $|s_2 - s_1|$. The difference $s_2 - s_1$ indicates both the distance and direction of s_2 from s_1 .

Example 2-4a. A certain motion is described, from the time $t = 5$ until the time $t = 8$, by the equation $s = \phi(t) = 2t^3 - 39t^2 + 252t - 535$.

The distance measurements s in a motion of this sort may be graphed on a single straight line, an s -axis, and for this reason we simplify our discussion by talking as though the motion actually took place on a straight line, although the path of the motion may make many turns, go up and down hills, etc. The moving object, which may, for example, be a large automobile, is represented as a single geometric point on the s -axis.

We need a concept similar to average velocity but valid for a precise instant $t = t_0$; by way of distinction, we call it the instantaneous velocity, or simply the velocity. Once again we have introduced a concept which cannot be defined exactly without methods of the calculus. We try to approximate the instantaneous velocity at time t_0 for a motion described by $s = \phi(t)$. We take the average velocity over a time interval between the time t_0 and another time t of the motion, that is

$$\frac{\phi(t) - \phi(t_0)}{t - t_0}$$

If these approximations have a limit, this limit is the velocity at time t_0 of the motion described by $s = \phi(t)$. That is, we shall define the velocity to be the number to which this ratio approximates within any prescribed error tolerance ϵ for values of t chosen sufficiently near to t_0 .

Example 2-4b. For the motion of Example 2-4a a calculation similar to that of Section 2-2 yields the velocity $6t^2 - 78t + 252$ (proof left as an exercise). At the time $t = 5.5$, for instance, the instantaneous velocity is 4.5 (distance units per time unit). This positive number cannot possibly be approximated within a small tolerance by negative average velocities; therefore over any sufficiently small time interval containing $t = 5.5$, the average velocity was positive, so that a distance in the positive direction of motion was covered during the time. It follows that the motion at $t = 5.5$ was in the forward direction.

We see that the velocity helps us to answer such questions as: when is the object moving in the direction of increasing s , when in the direction of decreasing s ? When is the object, at least momentarily, standing still (as it must do, for example, at an instant when it is reversing direction)?

To check for reversals of the motion, we factor the expression for the velocity: $6(t - 6)(t - 7)$. We discover that the velocity was zero at times $t = 6$ and $t = 7$, but at no other times during the motion. The corresponding locations are $s = \phi(6) = 5$ and $s = \phi(7) = 4$. This is possible only if the motion during the time interval $6 < t < 7$ was in the direction of decreasing s , and the velocity -1.5 obtained for $t = 6.5$ checks this assertion. The motion ended at time $t = 8$, and the final location was $s = \phi(8) = 9$. This is possible only if the point $s = 5$ was reached twice during the motion. (How does this conclusion follow?)

We tabulate our findings:

Time	5.0	5.5	6.0	6.5	7.0	7.5	8.0
Location			5.0		4.0		9.0
Velocity		4.5	0.0	-1.5	0.0		

As the example suggests, the velocity is positive at a time when the motion is in the positive direction and negative at a time when the motion is in the negative direction. At each instant when the direction of motion reverses, the velocity must necessarily be zero, but it may also be zero during a motion that never goes backward (see Exercises 2-4, No. 7 below).

The representation $s = \phi(t)$ of motion on a straight line has a somewhat different character from the representation $y = f(x)$ of the graph of a function f in a plane with coordinates x and y . In the representation of a graph, the coordinates x and y from the geometrical point of view assume equal importance, and are essentially equivalent. In the representation $s = \phi(t)$ of a motion, the quantity t assumes the role, not of a coordinate equivalent to s , but of a parameter, a secondary mathematical quantity in terms of which the primary quantity, the distance s , is expressed. Look carefully at the contrast between coordinate and parametric representation as employed on a single axis, an x -axis or s -axis. In the former, a point P is specified by a number x which directly pinpoints the location of P on the line. In the latter, a location Q is specified by a number t having no direct geometrical significance. The quantity of geometrical significance is the number $\phi(t)$, the value at t of the parametric function ϕ , and it reveals the location s . The advantage here is that a description such as that of Example 2-4a tells us not merely the location of the moving object, but also the specific time at which it occupied the stated location. The description $s = 5$ would not distinguish between the two times when the moving object passed through this location, but $t = 6$ precisely describes the first time, and $t = 7.5$ the second time.

Parametric representation has many other applications; the parameter may in another context represent distance, or temperature, or heat, or fluid speed. As will be seen in Chapter 11, the parametric representation of a curve in the plane may be far simpler than its coordinate representation in terms of x and y , and enables us to study motion on plane curves in a much more efficient way.

Exercises 2-4

1. Let us assume that a pellet is projected straight up and after a while comes straight down via the same vertical path to the place on the ground from which it was launched. After t seconds the pellet is s feet above the ground. Some of the ordered pairs (t, s) are given in the following table.

t	0	1	2	3	4	5	6	7	8	9	10
s	0	144	256	336	384	400	384	336			0

We shall intentionally avoid certain physical considerations such as air resistance. Moreover, we shall deal with simple numbers rather than quantities measured to some prescribed degree of accuracy which might arise from the data of an actual projectile problem in engineering.

- (a) Interpolate from the data given to determine the height of the projectile after eight and nine seconds respectively. (Guess, using symmetry as your guide.) Does extrapolation to find values of s for $t = -1$ or $t = 11$ make sense on physical grounds? After how many seconds does the projectile appear to have reached its maximum height? What seems to be the maximum height?
- (b) Does s appear to be a function of t ? If so, discuss the domain and range, taking physical considerations into account.
- (c) If we were to plot a graph of $s = f(t)$,
 - (1) is it plausible on physical grounds to restrict our graph to the first quadrant?
 - (2) Does the data suggest that the scale on the s -axis (vertical) should be the same as the scale on the t -axis (horizontal)?
- (d) Keeping in mind your responses to part (c), plot the ordered pairs (t, s) from the table. Connect the points with a smooth curve. What is the name of the function suggested by the graph? On physical grounds is it feasible that there would be a real value of s for every real number assigned to t over the interval $0 \leq t \leq 10$? Were we probably justified in connecting the points?
- (e) Assuming that the equation $s = f(t) = At^2 + Bt + C$ was used to develop the entries in our table, find values for constants A , B , and C .

- (f) Sketch the graph given by the equation $s = 160t - 16t^2$ over the interval $0 \leq t \leq 10$. Using a more carefully plotted graph of the above set, connect the point where $t = 1$ with the point where $t = 2$ with a chord. What is the slope of this chord? Estimate the slope of the curve at $t = 1$ and $t = 2$.
- (g) If the units of s are feet and the units of t are seconds, what are the units of slope? What word is commonly associated with this ratio of units? What would you guess are the physical interpretations of positive, zero, and negative values of this ratio?
- (h) Draw the graph of $v = 160 - 32t$ over the interval $0 \leq t \leq 10$. Compare the values of v for $t = 1$ and $t = 2$ respectively with your estimate for the slopes of the graph of $s = 160t - 16t^2$ in part (f).
- (i) Average the values of v for $t = 1$ and $t = 2$ and compare this average with the slope of the chord connecting the points where $t = 1$ and $t = 2$ in part (f).
- (j) If the units of v are ft./sec. and the units of t are seconds, what are the units of the slope of the line $v = 160 - 32t$? What word from physics is commonly associated with this ratio of units? Does the minus sign along with the particular numerical value of this slope have any special connotation from your experience?
2. (a) Derive the velocity function for the motion as given in Example 2-4a.
- (b) Sketch the graph of $s = \phi(t)$, (called the world line) and the v vs. t curve (i.e., the graph of the velocity as a function of time).
- (c) Compare the times when s equals a maximum or a minimum and when the velocity $v = 0$. Explain this physically.
- (d) Given only that $\phi(6) = 5$ for the function ϕ that describes the motion, show that there is a second time t when $\phi(t) = 5$, and find that value of t . (This is not done by calculus.)
- (e) Find the time of greatest speed between $t = 6$ and $t = 7$.

3. Find the velocity of an object whose location along a straight line is described by the equation $s = 128t - 16t^2$. Sketch the curves of s vs. t and v vs. t on the same set of axes.
- (a) During what time interval or intervals is the object moving toward the location $s = 0$?
- (b) What are the values of v and t when s is a maximum?
4. A ball is thrown upward with a velocity of 32 feet per second. Its height h in feet after t seconds is described by the equation $h = 32t - 16t^2$.
- (a) What is the velocity of the ball when its height first reaches 12 feet? When it again reaches 12 feet?
- (b) How high does it go, and at what time does the ball reach its highest position?
5. An object is projected up a smooth inclined plane in a straight line. Its distance s in feet from the starting point after t seconds is described by the equation $s = 64t - 8t^2$. After the object reaches its highest point it slides back along its original path to the starting point according to the equation $\bar{s} = 8(t - t_n)^2$. Here \bar{s} is the distance of the object from the highest point and t_n is the time it took the object to reach the highest point.
- (a) Determine how long it took for the object to make the up and down trip.
- (b) Sketch the s vs. t curve for the up and down motion using one set of coordinates. Do likewise for the v vs. t curve.
6. The location of an object on a straight line is given by the formula $s = pt^2 + qt + r$, where p , q , and r are real constants. Find all instants of time when the object is at rest, and show how the number of such instants depends on the constants p , q , and r .

7. For any but the very simplest motions, the function ϕ describing the location at time t is unlikely to be expressible in terms of a single formula for the entire duration of the motion. Here is a more plausible description of a motion:

$$\phi(t) = \begin{cases} \frac{t^2}{2}, & 0 \leq t \leq 1; \\ t - \frac{1}{2}, & 1 \leq t \leq \frac{5}{2}; \\ -\frac{t^2}{2} + \frac{7t}{2} - \frac{29}{8}, & \frac{5}{2} \leq t \leq \frac{7}{2}; \\ \frac{5}{2}, & \frac{7}{2} \leq t \leq 6; \\ t^2 - 12t + \frac{77}{2}, & 6 \leq t \leq 8; \\ \frac{21}{2} - 4(t - 7)^{-1}, & 8 \leq t \leq 10. \end{cases}$$

- (a) Compute the function $\psi(t)$ that describes the velocity of the motion with location $s = \phi(t)$.
- (b) It is claimed that $s = \phi(t)$ and $v = \psi(t)$ are functions. What has to be checked to verify this? Does it check? Show the graph of each of these functions on the same axes.
- (c) During what time intervals do you think the speed of the motion is increasing? decreasing?
- (d) Does the object spend any time between $t = 0$ and $t = 10$ standing still? Does it have any other instants of zero velocity during the motion?

2-5. Derivative of a Function at a Point.

Slope and velocity are two concepts with seemingly unrelated origins in geometry and physics; yet we have obtained each in the same way as a limit. This sort of limit is one of the most important ideas of mathematics, and also appears in countless guises in the sciences, in technology, and in the social sciences. We abstract from the ideas of velocity and slope a general analytical concept--the derivative of a function at a point--and frame this concept in a purely numerical way, independently of the problems from which it arises.

Let a be a value of x in the domain of a function f . We say that m is the derivative of f at a if, by taking x sufficiently near to a , we can approximate m by the ratio $\frac{f(x) - f(a)}{x - a}$ with an error less than any prescribed tolerance. We briefly describe this situation by saying that m is the limit of $\frac{f(x) - f(a)}{x - a}$ as x approaches a , and symbolically we write

$$m = \lim_{x \rightarrow a} \frac{f(x) - f(a)}{x - a}.$$

(Read: " m is the limit as x approaches a of $\frac{f(x) - f(a)}{x - a}$.")*

The procedure of computing a derivative is called differentiation, and if f has a derivative at a --if the limit exists--then we say that f is differentiable at a .

The preceding description gives no prescription for differentiating any particular function at any particular point. In fact, no universal technique for finding derivatives exists, although with the assistance of the general methods of operation with limits to be developed in the next chapter we shall be able to differentiate many of the functions of greatest interest. For the present we try to pursue a little further the technique which enables us to find the derivatives of simple polynomial functions in Sections 2-2 and 2-3.

Look again at the ratio $r(x) = \frac{f(x) - f(a)}{x - a}$. No matter what function f appears in the numerator, the denominator forbids us to evaluate r at $x = a$. In each previous example we have avoided this difficulty by finding

*The ratio $\frac{f(x) - f(a)}{x - a}$ is often thought of as the ratio of the "change" in the values to f to the "change" in the values of the independent variable, and is referred to as the "average rate of change" of $f(x)$. The limit m is then called the "rate of change" of $f(x)$ at $x = a$.

a simple formula for a function which had the same values as r for $x \neq a$, but which was also defined for $x = a$. Thus for $f(x) = x^2$ we found

$$r(x) = \frac{x^2 - a^2}{x - a} = \frac{(x - a)(x + a)}{(x - a)} (x + a)$$

or

$$r(x) = x + a, \quad \text{for } x \neq a.$$

This formula suggests that the limit of $r(x)$ for x approaching a is $m = 2a$.

In finding the value of m , the crucial step was the representation of the numerator of $r(x)$ as a product containing the factor $(x - a)$. With the help of algebraic techniques it is possible to find the derivatives of a number of simple functions in this way.

Example 2-5a. Consider the function g defined by $g(x) = \frac{1}{x}$, whose derivative we wish to find at $x = a$. As yet, we cannot even guarantee the existence of a derivative without further investigation (see, for instance, Exercises 2-3, No. 8). For the function g we have

$$r(x) = \frac{g(x) - g(a)}{x - a} = \frac{\frac{1}{x} - \frac{1}{a}}{x - a} = \frac{-x + a}{ax(x - a)} = -\frac{1}{ax} \frac{(x - a)}{(x - a)},$$

or

$$r(x) = -\frac{1}{ax}, \quad \text{for } x \neq a.$$

Our natural guess for the derivative at a is $m = -\frac{1}{a^2}$. To test our guess we should set a tolerance $\epsilon > 0$, and check that the error in the approximation of m by $r(x)$ remains within that tolerance for x sufficiently close to a . That is, we should show that for such x we have

$$\left| r(x) - m \right| = \left| r(x) + \frac{1}{a^2} \right| = \left| \frac{x - a}{a^2 x} \right| = \frac{|x - a|}{a^2 |x|} < \epsilon, \quad (x \neq a).$$

The situation is not quite so simple as before. The error is still a multiple of $|x - a|$, but the multiplier depends on both a and x ; and we must keep x away from 0. Algebraic techniques for completing the verification are developed in Section 3-3. You are encouraged to try to complete it yourself (Exercises 2-5, No. 9a). We repeat our result: the derivative at $x = a$ of $g : x \rightarrow \frac{1}{x}$ is $m = -\frac{1}{a^2}$.

Example 2-5b. Another function whose derivative at a can be obtained with the aid of elementary algebra is the square root function

$h : x \rightarrow x^{\frac{1}{2}} = \sqrt{x}$. For our approximations to the derivative m we find now

$$r(x) = \frac{h(x) - h(a)}{x - a} = \frac{\sqrt{x} - \sqrt{a}}{x - a}, \quad (x \neq a).$$

In order to obtain the factor $(x - a)$ in the numerator, we rationalize the numerator.

$$r(x) = \frac{(\sqrt{x} - \sqrt{a})(\sqrt{x} + \sqrt{a})}{(x - a)(\sqrt{x} + \sqrt{a})} = \frac{(x - a)}{(x - a)} \frac{1}{\sqrt{x} + \sqrt{a}}, \quad (x \neq a).$$

The value of $r(x)$ is $\frac{1}{\sqrt{x} + \sqrt{a}}$ for $x \neq a$, and $m = \frac{1}{2\sqrt{a}}$ is our guess for the derivative of h at $x = a$. The proof that this is the limit is similar to the verification for the preceding function, and is also left for later (see Exercises 2-5, No. 9b and Section 3-3). For the present we accept the result that the derivative at $x = a$ of $h : x \rightarrow \sqrt{x}$ is $\frac{1}{2\sqrt{a}}$.

Example 2-5c. Before you reach the wrong conclusion that all derivatives may be obtained with nothing more than simple algebra, take a careful look at the sine function, $k : x \rightarrow \sin x$. While the domain of the sine function consists of all real x , for simplicity we shall attempt to compute the derivative only at $x = a = 0$. We then need to find the limit as x approaches 0 of the ratio

$$r(x) = \frac{k(x) - k(0)}{x - 0} = \frac{\sin x - \sin 0}{x - 0} = \frac{\sin x}{x}.$$

How can we divide the trigonometric function $\sin x$ by the linear polynomial x ? The answer, unfortunately for our attempt to differentiate $\sin x$, is that there is no method of algebra or trigonometry that enables us to carry out such a division. No trick will help, and we are forced to take a quite different route in order to find the derivative. After a more detailed study of limits, and the application of this study to a more systematic calculation of the derivatives of algebraic functions, we shall return to this problem and solve it.

Exercises 2-5

1. Find the derivative of f at a for $f : x \rightarrow f(x)$, where $f(x)$ equals each of the following.
 - (a) $\frac{1}{x+1}$
 - (b) $\frac{1}{Ax+B}$, where A and B are constants, $A \neq 0$
 - (c) $\frac{1}{x^2}$
 - (d) $\frac{1}{x^3}$
 - (e) $\frac{1}{\sqrt{x}}$
 - (f) $x^{\frac{3}{2}}$

2. For each of the functions of Exercise 1, and also for the functions $g(x) = \frac{1}{x}$ and $h(x) = \sqrt{x}$ of the text, list the following: the domain of the function; all points of horizontal slope (if any); the highest point on the graph of the function (if any); the lowest point on the graph of the function (if any).

3. Use the definition of derivative to differentiate f at both $a = -2$ and $a = 2$ if $f(x)$ equals
 - (a) $(x-2)^2$
 - (b) $\frac{1+x}{1-x}$
 - (c) $3x - x^3$
 - (d) $x^2 - x - x^{-1}$

4. For each of the functions f whose values $f(x)$ are described below, find the derivative at a , where a is in the domain of the function.

(a) $5x - x^2$	(d) $\frac{2x}{x+1}$
(b) $3 - x - x^2$	(e) $\frac{2}{x-1}$
(c) $x^3 - 2x$	(f) $\frac{x^2 + x}{x^2 - 1}$

5. By the methods of this section find the slope of the graph of each of the following functions at the point $(1,1)$.

(a) $f : x \rightarrow x^3$

(b) $g : x \rightarrow x^{\frac{1}{3}}$

What is the relationship between the two answers? Explain this relationship.

6. Use a table of sines to obtain approximations to the limit of $\frac{\sin x}{x}$ as x approaches $a = 0$, and make a conjecture as to the value of this limit. What can you conclude about the slope of the graph of $y = \frac{\sin x}{x}$ at the origin?
7. Write the derivative of the cosine function at $x = a$ as a limit, for $a = 0$ and $a = \frac{\pi}{2}$, and show what limits must be evaluated before the derivative can be obtained.
8. In Chapter 8 we shall define a function $f : x \rightarrow 2^x$ in such a way as to have all real numbers x in its domain, and also to satisfy all the familiar laws of exponents valid for 2^n when n is a rational number. Compute the ratios $r(x)$ that serve to approximate the derivative of f at $x = a = 0$, and show what limit problem must be solved before we can differentiate the function.
9. Prove that the derivative at $x = a$ of
- (a) $g : x \rightarrow \frac{1}{x}$ is $-\frac{1}{a^2}$
- (b) $h : x \rightarrow \sqrt{x}$ is $\frac{1}{2\sqrt{a}}$.
10. Find the slopes (if they exist) of the following curves at points for which $x = y$:
- (a) $x + y = A$
- (b) $x + xy + y = 3A$
- (c) $xy = A^2$
- (d) $|x + y| + |x - y| = 2A$
- (e) Check your answers by sketching each of the above curves.
- (f) Generalize the results of parts a - d.

Chapter 3

LIMITS AND CONTINUITY

3-1. Introduction.

The derivative is a special example of a basic general concept: the limit of a function at a point. The concept of limit is also closely related to the important idea of continuous function which is taken up in the later sections of this chapter.

In simple cases, a judicious guess may help us find limits without apparent need for deductive reasoning. As we widen our mathematical horizons and try to solve more complicated problems with the calculus, we soon exceed the potentialities of such methods. Only a firmly founded theory of limits will enable us fully to exploit the methods of the calculus. It is our purpose to provide such a foundation in this chapter.

It is one of the triumphs of the calculus that an enormous variety of significant problems can be solved by straightforward formal operations which readily yield solutions in terms of derivatives or integrals. In the formal calculus the idea of limit and the process of approximation upon which it is based disappear. Of course, there are always problems for which the formal techniques fail; for these the underlying theory becomes essential and we must go back to first principles.

Naturally you cannot expect to develop the theory of limits and thus acquire a deeper understanding of the calculus without an expenditure of hard thinking and careful work. Moreover, it is unusual to obtain a complete grasp of such a subtle concept at first encounter. The study of mathematics is a slow maturing process, and you can be well satisfied if you comprehend enough of the theory in this chapter to enable you to follow its later applications. As you grow in mathematical skill and maturity, a more complete understanding will come; for the present, try to see how we are able to use the theory for the practical purpose of developing a calculus--a scheme of rapid and efficient reckoning--of limits, derivatives, and later, integrals.

The most important use of the basic theory is to approximate solutions of complicated problems in terms of simple functions. Such approximations are of special interest in this day of high-speed computers. Although a problem may have a complete formal solution by the methods of the calculus, for the purpose of numerical computation it may save time, effort, and money to ignore the formal solution and treat the problem by approximation. Simple approximations play a fundamental role in science and engineering. A realistic model of a phenomenon may involve so many complexities that the problem is intractable mathematically. As we shall see in later chapters, a simplified model yielding an approximation to the complete solution is often more useful. Limits enter in approximations to the solution of the simplified problem when there is no formal explicit solution, and in estimation of the error when the unobtainable solution of the realistic problem is replaced by the solution of the simplified problem.

Exercises 3-1

1. Consider $f : x \longrightarrow [x] + [-x]$. (See Section A2-1 for discussion of $[x]$.)
 - (a) What number if any do the values of f approximate when x is close to 1? When x is close to 2?
 - (b) What can you say of the limit of $f(x)$ as x approaches n , when n is any integer?
 - (c) Evaluate $f(n)$, where n is an integer.
 - (d) Sketch the graph of f . Go back and check your answers to (a) - (c). Do they agree with your graph?
2. For each of the following functions sketch the graph and, if possible, find the limit as x approaches 0.
 - (a) $f : x \longrightarrow \frac{x}{x}$
 - (b) $f : x \longrightarrow \frac{1}{|x|}$
 - (c) $f : x \longrightarrow \frac{x}{|x|}$
 - (d) $f : x \longrightarrow \frac{[x]}{x}$
 - (e) $f : x \longrightarrow \sin \frac{1}{x}$

3-2. Definition of Limit of a Function

In Chapter 2 we used the concept of limit of a function in defining derivative. At this point we need a precise formulation of the limit concept in order to derive the laws which govern operations with limits.

Although the concept of limit of a function is more general than the idea of derivative, our study of limits was initially motivated by the basic example of the derivative of a function ϕ as the limit of the ratio $r(x)$:

$$m = \lim_{x \rightarrow a} r(x),$$

where

$$r(x) = \frac{\phi(x) - \phi(a)}{x - a}.$$

In order to be sure that the description of the derivative as the limit of $r(x)$ makes sense we must be sure that we have an adequate set of approximations, that $r(x)$ is defined for numbers x arbitrarily close to a .

Usually, the domain of r will contain an entire neighborhood of a (excluding a itself) but either for theoretical or practical reasons it is often useful to analyze the behavior of $r(x)$ on only one side of a . For example, there is a natural starting point in the motion of a rocket and it is essential to know the initial direction of the rocket in order to determine the rest of the trajectory.*

In framing the general definition of the limit of a function f at a point a we then require that we have an adequate set of approximations. Specifically, the definition may not include the value $f(a)$ among the approximations, even if it should be defined, but it must involve values $f(x)$ for x close to a . For this purpose we introduce the deleted h -neighborhood of a , that is, the set of all x for which

$$0 < |x - a| < h.$$

As the set of approximations to be used in defining the limit of f at a we take the set of values $f(x)$ for all x from the domain of f in some deleted neighborhood of a .

*In some texts this important case is taken care of by separate definitions of "right-sided" and "left-sided" limits. (See Exercises 3-4, No. 16.)

With these ideas in mind we are now able to express the idea of limit completely in analytical terms. If f has a limit L as x approaches a then for any error tolerance ϵ we keep $f(x)$ within ϵ of L by restricting x to be any number from the domain of f in a sufficiently small neighborhood of a .

DEFINITION 3-2.* Let a be a point for which every deleted neighborhood contains points of the domain of f . The function f has the limit L at a if (and only if) for each positive number ϵ , there exists a positive number δ such that

$$|f(x) - L| < \epsilon$$

for every x in the domain of f which satisfies the inequality

$$0 < |x - a| < \delta.$$

We then write $\lim_{x \rightarrow a} f(x) = L$.

It follows from the definition of limit, since the value $f(a)$ itself does not lie in the class of approximations considered, that any function which takes on the same values as f in some deleted neighborhood of a would have the same limit at a . For example, the two functions f and g defined below have the same limit at every point a of the real axis.

$$f(x) = 1$$

$$g(x) = \begin{cases} 0, & \text{for an integer } x, \\ 1, & \text{for non-integral } x. \end{cases}$$

Although we do not rely upon pictures for our precise understanding of the concept of limit, it is desirable to have a geometrical interpretation of the idea.

Example 3-2a. The graph of the function

$$f: x \rightarrow 2x - 4$$

*The definition of limit can be recapitulated in terms of neighborhoods: the number L is said to be the limit of f at a if every deleted neighborhood of a contains points of the domain of f and if for each ϵ -neighborhood of L there is at least one deleted δ -neighborhood wherein f maps the points of its domain into the ϵ -neighborhood.

is shown in Figure 3-2a. In order to show that

$$\lim_{x \rightarrow 3} (2x - 4) = 2$$

we must show, for every $\epsilon > 0$, that there is a $\delta > 0$ so that

$$|(2x - 4) - 2| < \epsilon$$

for all x in the deleted neighborhood $0 < |x - 3| < \delta$. It is easy to see from Figure 3-2a how δ may be found.

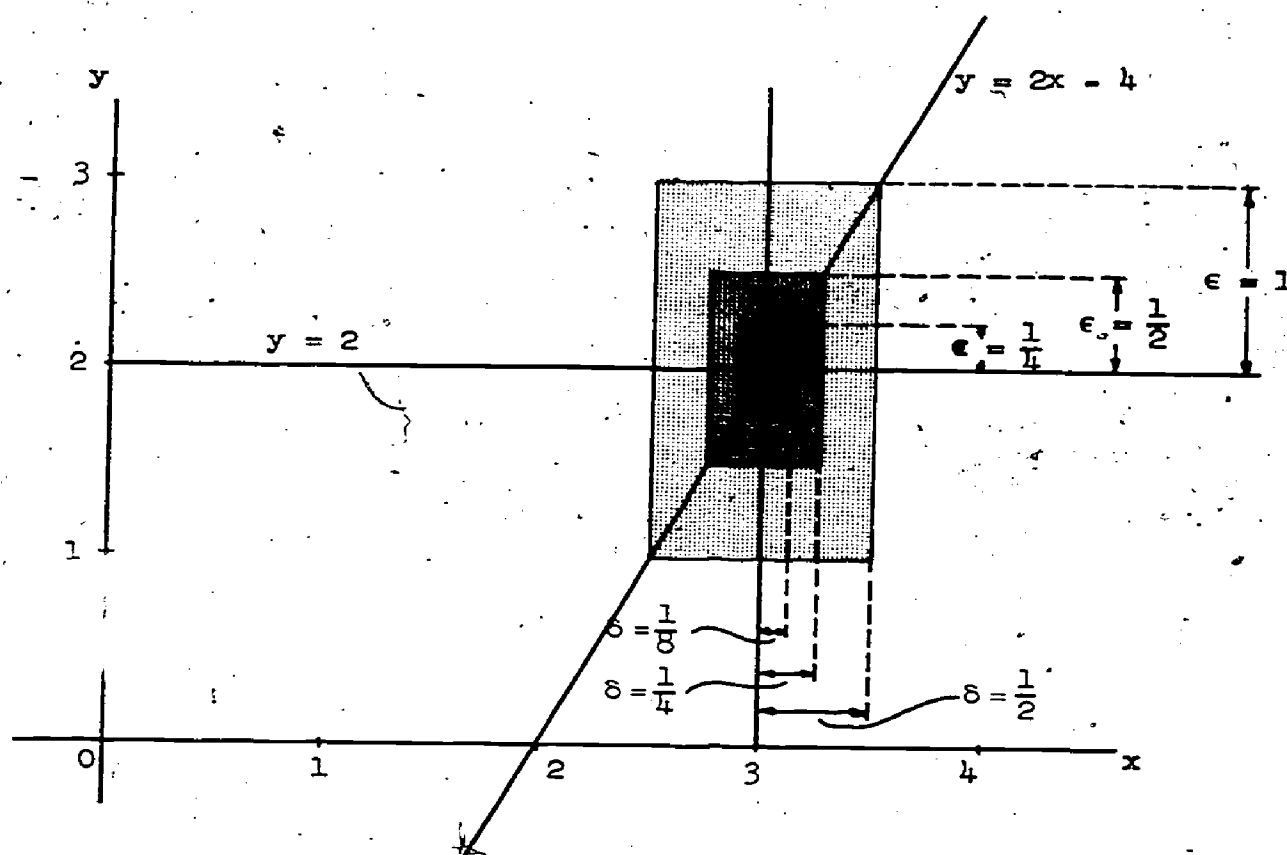


Figure 3-2a

Given a horizontal band of width 2ϵ centered on the line $y = 2$, we can find a vertical band of width 2δ about $x = 3$ so that the graph of f lies entirely within the rectangle where the bands overlap. From the graph we infer that for $\epsilon = 1$ we may take $\delta = \frac{1}{2}$, for $\epsilon = \frac{1}{2}$, $\delta = \frac{1}{4}$, and for $\epsilon = \frac{1}{8}$, $\delta = \frac{1}{8}$. There seems to be no obstacle to finding a δ for any ϵ , no matter how small, but we clearly cannot rely on pictures to do so. Instead, we proceed analytically. If we require $0 < |x - 3| < \delta$, then

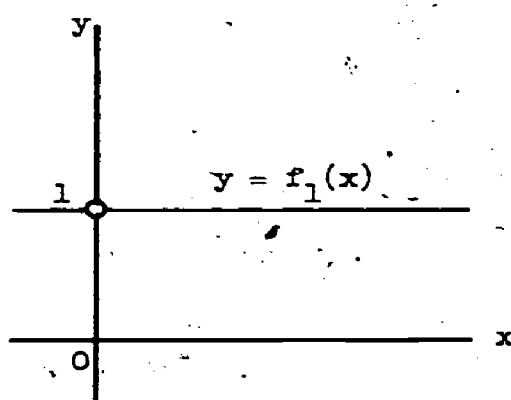
$$\begin{aligned}
 |f(x) - 2| &= |(2x - 4) - 2| \\
 &= |2x - 6| \\
 &= |2(x - 3)| \\
 &= 2|x - 3| \\
 &< 2\epsilon .
 \end{aligned}$$

Consequently, if we take $\delta = \frac{\epsilon}{2}$, then

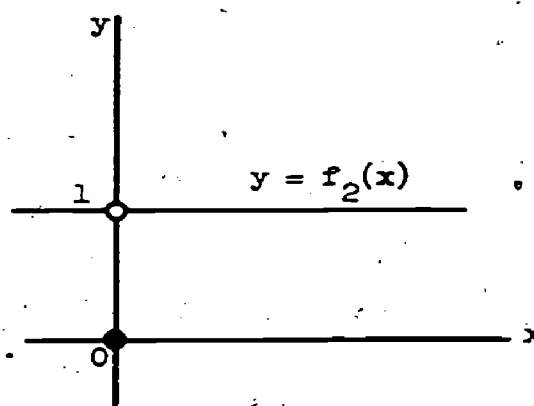
$$|f(x) - 2| < \epsilon .$$

The preceding example was made especially simple to reveal the basic picture. We now explore the concept of limit in a variety of situations.

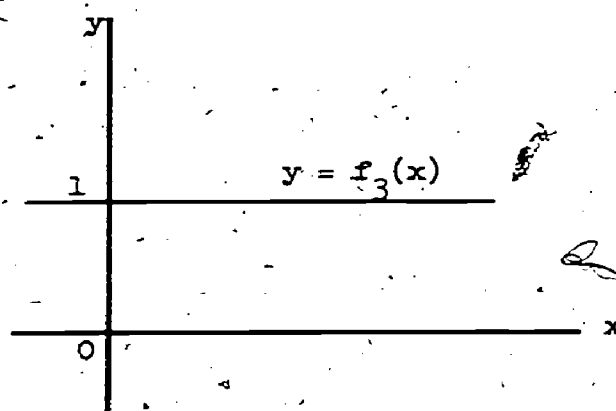
Example 3-2b. Figure 3-2b presents the graphs of the three functions given by $f_1(x) = \operatorname{sgn}(\frac{1}{x^2})$, $f_2(x) = \operatorname{sgn} x^2$, $f_3(x) = 1$.



$$f_1 : x \rightarrow \operatorname{sgn}\left(\frac{1}{x^2}\right)$$



$$f_2 : x \rightarrow \operatorname{sgn} x^2$$



$$f_3 : x \rightarrow 1$$

Figure 3-2b

Observe that $x = 0$ is a point of the domains of f_2 and f_3 but not of f_1 . For each of these functions we wish to consider the limit, if it exists, as x approaches 0.

Since the three functions coincide when $x \neq 0$, and the value of the limits does not depend on how the functions are defined at $x = 0$, it is clear that all three functions have the same limit. In each case 1 is the obvious candidate for the limit. Verify that the conditions of Definition 3-2 are satisfied by $L = 1$ at $x = 0$.

Observe that there is a gap in the graphs of f_1 and f_2 at $x = 0$, and that the graph of f_3 is continuous, it has no gap. The function f_1 has a limit at $x = 0$, but is not defined there, f_2 is defined at $x = 0$ but $f_2(0)$ is not its limit, f_3 has a limit at $x = 0$ and the limit is the function value. We see then that the concept of limit and the intuitive idea of continuity are closely related; we shall pursue this connection further in Section 3-5.

Example 3-2c. Figure 3-2c presents the graphs of the two functions given by

$$g_1(x) = x^2 + \operatorname{sgn}(x - a)$$

$$g_2(x) = x^2 + \operatorname{sgn} \sqrt{x - a}$$

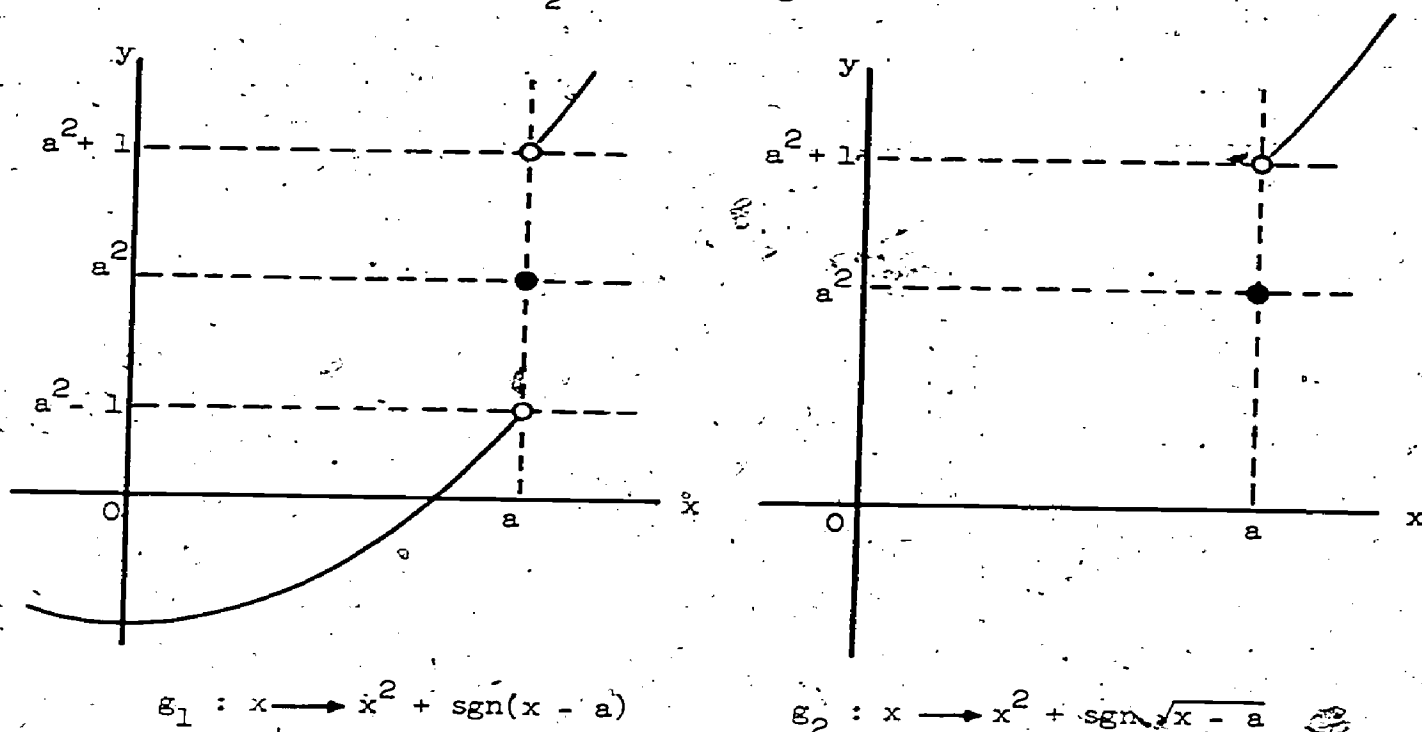


Figure 3-2c

The function g_1 is defined for all values of x . The domain of g_2 consists only of those values of x for which $x \geq a$ and on this domain it has the same values as g_1 . It seems clear from the graph that there is no single number L which is approximated by the values $g_1(x)$ as x approaches a . On the contrary, in any neighborhood of a , it is possible to find values of x for which $g_1(x)$ approximates $a^2 - 1$ within any given error tolerance and other values which approximate $a^2 + 1$. Verify, then, that the conditions of Definition 3-2 cannot be satisfied, that g_1 has no limit at $x = a$.

For the function g_2 , on the other hand, it appears that no matter what the error tolerance, there is a deleted neighborhood of a wherein $g_2(x)$ approximates $a^2 + 1$ within the tolerance for all x in the domain of the function. This is easily verified. In a deleted δ -neighborhood of a , we have

$$g_2(x) = x^2 + 1, \text{ for } a < x < a + \delta.$$

We have for the absolute error of approximation

$$\begin{aligned} |g_2(x) - (a^2 + 1)| &= |x^2 - a^2| \\ &= |x - a| \cdot |x + a| \\ &< \delta(|x| + |a|) \\ &\leq \delta((|a| + \delta) + |a|) \\ &< \delta(2|a| + \delta). \end{aligned}$$

This absolute error can be kept within any given error tolerance ϵ by restricting x to a small enough δ -neighborhood of a . For simplicity, we first restrict ourselves to neighborhoods of radius no larger than 1. Taking $\delta \leq 1$ in the inequality above, we obtain a simpler bound on the absolute error in terms of the radius δ :

$$|g_2(x) - (a^2 + 1)| < \delta(2|a| + 1).$$

Now if we choose δ so that

$$\delta \leq \frac{\epsilon}{2|a| + 1},$$

then we have ensured that

$$|g_2(x) - (a^2 + 1)| < \epsilon,$$

namely, that the error has been kept within the tolerance ϵ . Since this is a prescription for controlling the error within any tolerance ϵ , we have accomplished our purpose and proved

$$\lim_{x \rightarrow a} g_2(x) = a^2 + 1$$

completely in the analytic terms of Definition 3-2.

Exercises 3-2

1. Show that if $0 < |x - a| < 1$, then $|x + 2a| < 1 + 3|a|$.
2. Show that if $0 < |x - a| < 1$, then $|x^3 - a^3| < (3|a|^2 + 3|a| + 1)|x - a|$.
3. Show that if $0 < |x - 2| < 1$, then $\frac{1}{|x - 4|} < 1$.

Hint: By Exercise 10 of A1-2, if $|x - 4| > 1$, then $\frac{1}{|x - 4|} < 1$.

4. Show that if $|x - a| < \frac{|a|}{2}$, then $\frac{1}{x^2} < \frac{4}{a^2}$.
5. Show that if $0 < |x - 1| < 1$, then $|4x + 1| < 9$ and $\left| \frac{1}{x + 2} \right| < 1$.
6. Show that if $0 < |x - 2| < 1$, then $|x + 1| < 4$ and $\left| \frac{1}{x^2 + 2x + 4} \right| < 1$.
7. Estimate how large $x^2 + 1$ can become if x is restricted to the open interval $-3 < x < 1$.
8. Use inequality properties to find a positive number M such that $0 < |x - 1| < 3$ for all x and,
 - (a) $|x^2 + 2x + 4| \leq M$.
 - (b) $|3x^2 - 2x + 3| \leq M$.
9. (a) Show that if $0 < |x - 3| < 1$ and $0 < |x - 3| < \frac{\epsilon}{7}$, then $|x^2 - 9| < \epsilon$.
 - (b) Show that the pair of inequalities $\delta \leq 1$ and $\delta \leq \frac{\epsilon}{7}$ (or $\delta \leq \min\{1, \frac{\epsilon}{7}\}$) is satisfied by $\delta = \frac{\epsilon}{7 + \epsilon}$.
10. Find a number $M \geq 1$ such that $\left| \frac{x + 4}{x - 4} \right| \leq M$ for all x such that $0 < |x - 2| < 1$. (See No. 3.)
11. For the given value of ϵ , find a number δ such that if $0 < |x - 3| < \delta$, $|x^2 - 9| < \epsilon$.
 - (a) $\epsilon = 0.1$.
 - (b) $\epsilon = 0.01$.

Is your choice of δ in (b) acceptable as an answer in (a)? Explain.

12. For the following functions, find the limit L as x approaches a .
For each value of ϵ , exhibit a number δ such that $|f(x) - L| < \epsilon$
whenever $|x - a| < \delta$.

(a) $f(x) = 3x - 2$, $a = \frac{1}{2}$.

(b) $f(x) = mx + b$, ($m \neq 0$)

(c) $f(x) = 1 + x^2$, $a = 0$

3-3. Epsilonic Technique.

It is conventional in discussions of approximations to a limit to use the Greek letter epsilon for the error tolerance. For this reason the subject devoted to techniques for the control of error is colloquially called epsilonics. We shall make immediate use of epsilonic technique in deriving the limit theorems which follow this section. Eventually, in applications, skill in epsilonic technique will be extremely valuable for making estimates when it is difficult to work with precise values. To develop this skill it is helpful to set up a routine pattern in which to present an epsilonic argument. We shall first describe the pattern in general and then, for several examples, carry out the proofs as indicated in the pattern.

Statement of the problem.

To prove that $\lim_{x \rightarrow a} f(x) = L$:

For each tolerance: $\epsilon > 0$ obtain a control δ .

Show: if $0 < |x - a| < \delta$, then $|f(x) - L| < \epsilon$.

We have stated the problem in outline. The proof is based on Definition 3-2. We must control the error $|f(x) - L|$ within the error tolerance ϵ by restricting the values of x to a sufficiently small deleted neighborhood of a . The proof is completed by verifying for a suitable radius δ that it does give the desired degree of control. The crucial open question is, how do we choose a suitable δ ?

Step 1. Simplification.

Find a $g(\delta)$ such that if $0 < |x - a| < \delta$, then $|f(x) - L| < g(\delta)$.

The idea here is to obtain an upper bound $g(\delta)$ for the absolute error where $g(\delta)$ can be held within the tolerance ϵ by taking sufficiently small values of δ . If we have $g(\delta) \leq \epsilon$, then $|f(x) - L| < g(\delta) \leq \epsilon$ and our objective is achieved. In some of the following examples the work of simplification is divided into three stages: (a) $f(x)$ is expressed in terms of $x - a$; (b) from the inequality $0 < |x - a| < \delta$ there is derived an inequality of the form $|f(x) - L| < g(\delta)$; (c) a δ is chosen for each ϵ in such a way that $g(\delta) \leq \epsilon$. In general, g is to be a simple function, one for which it is easy to find a δ such that $g(\delta) \leq \epsilon$. More typically, it will even be possible to solve for δ in the equation $g(\delta) = \epsilon$. For most of the cases in this text it is possible to obtain $g(\delta) = c\delta$ with a positive constant of proportionality c . Manipulations yielding a simple expression for $g(\delta)$ are illustrated in the Examples below.

Step 2. Choice of δ :

Choose δ so that $g(\delta) \leq \epsilon$.

This is the place where the work of simplification in Step 1 pays off. In

the most typical case where $g(\delta) = c\delta$ we may choose $\delta = \frac{\epsilon}{c}$.

Steps 1 and 2 show how the solution is found. The next step is the actual proof where we verify that the solution has been found.

Step 3. Verification.

Return to the statement of the problem. From the given expression for δ deduce the conclusion.

First we try out the method in a case where no complications arise, the case of the general linear function.

Example 3-3a.

Statement of the problem.

To prove that $\lim_{x \rightarrow a} (mx + b) = ma + b$, $(m \neq 0)$.

For each $\epsilon > 0$ obtain a δ .

Show: if $0 < |x - a| < \delta$, then $|f(x) - L| < \epsilon$.

Step 1. Simplification.

$$\begin{aligned} (a) \quad f(x) - L &= (mx + b) - (ma + b) \\ &= m(x - a). \end{aligned}$$

$$(b) \quad \text{If } |x - a| < \delta,$$

$$\begin{aligned} |f(x) - L| &= |m(x - a)| \\ &= |m| \cdot |x - a| \\ &< |m| \delta. \end{aligned}$$

$$(c) \quad \text{Take } g(\delta) = |m|\delta.$$

Step 2. Obtain δ .

To make $g(\delta) \leq \epsilon$, set

$$\delta = \frac{\epsilon}{|m|}$$

(allowable, since $|m| \neq 0$ by assumption).

Step 3. Verification.

Enter the result $\delta = \frac{\epsilon}{|m|}$ in the statement of the problem. The verification follows the pattern of Step 1 with one additional step:

$$\begin{aligned}
 |f(x) - L| &< |m|\delta \\
 &< |m| \frac{\epsilon}{|m|} \\
 &< \epsilon
 \end{aligned}$$

Since there is a strong inequality in this chain, we have

$$|f(x) - L| < \epsilon$$

In the following examples we shall omit repetitious material.

Example 3-3b.

Statement of the Problem.

To prove that $\lim_{x \rightarrow 0} \frac{1}{1 + |x|} = 1$.

For each $\epsilon > 0$ obtain a δ .

Show: if $0 < |x - 0| < \delta$, then $\left| \frac{1}{1 + |x|} - 1 \right| < \epsilon$.

Step 1.

$$\begin{aligned}
 \text{(a)} \quad \frac{1}{1 + |x|} - 1 &= \frac{1}{1 + |x|} - \frac{1 + |x|}{1 + |x|} \\
 &= \frac{-|x|}{1 + |x|}
 \end{aligned}$$

(b) If $0 < |x - 0| < \delta$,

$$\begin{aligned}
 \left| \frac{1}{1 + |x|} - 1 \right| &= \left| \frac{-|x|}{1 + |x|} \right| \\
 &= \frac{|x|}{1 + |x|} \\
 &< |x| \quad (\text{since } 1 + |x| > 1) \\
 &< \delta
 \end{aligned}$$

(c) Take $g(\delta) = \delta$.

Step 2. To make $g(\delta) \leq \epsilon$, set $\delta = \epsilon$.

Step 3. Set $\delta = \epsilon$ in the statement of the problem. We carry out the verification following Step 1 where we set $\delta = \epsilon$ at the last line.

The next example shows that it is not always sufficient to choose δ proportional to ϵ .

Example 3-3c.

Statement of the Problem.

To prove that $\lim_{x \rightarrow a} \sqrt{x} = \sqrt{a}$,

($a \geq 0$).

For each $\epsilon > 0$ obtain a δ .

Show: if $0 < |x - a| < \delta$ then $|\sqrt{x} - \sqrt{a}| < \epsilon$.

The choice $\delta = c\epsilon$, where c is a positive constant, cannot work when $a = 0$. In that case we observe that if $0 < x < \delta = c\epsilon$, then $\sqrt{x} < \sqrt{c} \sqrt{\epsilon}$.

We must then make $\sqrt{c} \sqrt{\epsilon} < \epsilon$ for all ϵ , no matter how small.

It follows that we must find a positive number c satisfying $\sqrt{c} < \sqrt{\epsilon}$ or, equivalently, $c < \epsilon$ for all positive ϵ . No such number exists; hence, $\delta = c\epsilon$ cannot work.

Step 1. From

$$|\sqrt{x} - \sqrt{a}| \leq |\sqrt{x} + \sqrt{a}| \quad (\text{Section A1-3, Formula 3})$$

we obtain on multiplying by $|\sqrt{x} - \sqrt{a}|$,

$$|\sqrt{x} - \sqrt{a}|^2 \leq |x - a|,$$

whence

$$|\sqrt{x} - \sqrt{a}| \leq \sqrt{|x - a|}.$$

Thus, if $0 < |x - a| < \delta$, then

$$|\sqrt{x} - \sqrt{a}| < \sqrt{\delta}.$$

Step 2. Choose $\delta = \epsilon^2$.

Step 3. Take $\delta = \epsilon^2$ in the statement of the problem. The verification is a recapitulation of Step 1 for this choice of δ .

It is often expedient to restrict δ by an auxiliary condition in Step 1. The following examples are typical.

Example 3-3d.

Statement of the Problem.

To prove that $\lim_{x \rightarrow -2} (x^3 - 5x - 1) = -3$.

For each $\epsilon > 0$ obtain a δ .

Show: if $0 < |x - 2| < \delta$, then $|(x^3 - 5x - 1) - (-3)| < \epsilon$.

Step 1.

$$(a) \quad x^3 - 5x - 1 - (-3) = x^3 - 5x + 2$$

$$= [(x - 2) + 2]^3 - 5[(x - 2) + 2] + 2$$

$$= (x - 2)^3 + 6(x - 2)^2 + 7(x - 2)$$

$$(b) \quad |x^3 - 5x - 1 - (-3)| = |(x - 2)^3 + 6(x - 2)^2 + 7(x - 2)|$$

$$= |(x - 2)\{(x - 2)^2 + 6(x - 2) + 7\}|$$

$$= |x - 2| \cdot |(x - 2)^2 + 6(x - 2) + 7|$$

$$\leq |x - 2| \cdot (|x - 2|^2 + 6|x - 2| + 7)$$

$$< \delta(\delta^2 + 6\delta + 7)$$

(At the last line we used $|x - 2| < \delta$.)

(c) For convenience we restrict δ by requiring $\delta \leq 1$. Under this condition

$$|x^3 - 5x - 1 - (-3)| < \delta(\delta^2 + 6\delta + 7)$$

$$\leq \delta(1 + 6 + 7)$$

$$\leq 14\delta$$

In order to get an upper bound in the simple form $c\delta$, we put a constant bound on the second factor in $\delta(\delta^2 + 6\delta + 7)$ by restricting δ . (The particular value 1 in $\delta \leq 1$ is inessential. We could have required $\delta \leq K$ where K is any positive constant.)

Step 2. We now wish to obtain a value δ satisfying two conditions simultaneously: $\delta \leq \frac{\epsilon}{14}$ and $\delta \leq 1$. One way of satisfying these conditions is to set

$$\delta = \frac{\epsilon}{14 + \epsilon}$$

where we have chosen the denominator simply as a convenient value which is greater than either 14 or ϵ . (See Exercise A1-3, No. 6a, b.)

Step 3. Set $\delta = \frac{\epsilon}{14 + \epsilon}$ in the statement of the problem. The verification follows Step 1 through (b). In (c) we use $\delta \leq 1$ and

$$\delta \leq \frac{\epsilon}{14}$$

to obtain

$$|(x^3 - 5x - 1) - (-3)| < \epsilon.$$

Alternative Step 1.

$$(a) \quad x^3 - 5x - 1 - (-3) = (x - 2)(x^2 + 2x - 1)$$

$$(b) \quad |x^3 - 5x + 2| = |x - 2| \cdot |x^2 + 2x - 1|$$

$$< \delta |x^2 + 2x - 1|$$

$$\leq 14\delta,$$

where, at the last line, imposing the condition $\delta \leq 1$ we utilize the result $1 < x < 3$ obtained from $|x - 2| < \delta \leq 1$.

Alternative Step 2. Since we do not use the formula for δ in the verification above but only the conditions $\delta \leq 1$ and $\delta \leq \frac{\epsilon}{14}$, it is natural (A1-3) to set

$$\delta = \min\left\{\frac{\epsilon}{14}, 1\right\}.$$

Alternative Step 3. Set $\delta = \min\left\{\frac{\epsilon}{14}, 1\right\}$ in the statement of the problem. The verification follows alternative Step 1 above.

From the preceding example we see that we have great freedom in choosing our control δ . We can always use more stringent controls than necessary: that is, given any deleted neighborhood of $|x - a| < \delta$, so chosen that $|f(x) - L| < \epsilon$ for any x in the neighborhood, then for all x in any subset of the neighborhood and, in particular, for any smaller deleted neighborhood of a , we satisfy the same inequality. In other terms, given any δ which keeps the error within the specified tolerance, any smaller value of δ will certainly have the same effect. It follows that we may impose the condition $\delta \leq K$ where K is any convenient positive constant. Similarly, having found a δ for a particular ϵ , we know that the same δ will suffice for any larger ϵ . Hence we need concern ourselves only with those ϵ satisfying $\epsilon \leq M$, where M is any convenient positive constant.

We conclude the list of examples by applying the techniques of the outline to prove that the conjectured values of the derivatives in Example 2-5a and Example 2-5b are, in fact, the correct limits.

Example 3-3e.Statement of the Problem.

To prove for $a \neq 0$ that

$$\lim_{x \rightarrow a} r(x) = \lim_{x \rightarrow a} \frac{\frac{1}{x} - \frac{1}{a}}{x - a} = \frac{1}{a^2} = L.$$

For each $\epsilon > 0$ obtain a δ .

Show: if $0 < |x - a| < \delta$, then $|r(x) - L| < \epsilon$.

(Observe that $r(x)$ is not defined at $x = 0$ or $x = a$.)

Step 1.

$$\begin{aligned} \text{(a)} \quad r(x) - L &= \frac{\frac{1}{x} - \frac{1}{a}}{x - a} + \frac{1}{a^2} \\ &= -\frac{1}{ax} \cdot \frac{x - a}{x - a} + \frac{1}{a^2} \\ &= -\frac{1}{ax} + \frac{1}{a^2}, \quad (x \neq a) \\ &= \frac{x - a}{a^2 x} \end{aligned}$$

(Note that we used $|x - a| > 0$ in setting $\frac{(x - a)}{(x - a)} = 1$ for $x \neq a$.)

$$\begin{aligned} \text{(b)} \quad |r(x) - L| &= \left| \frac{x - a}{a^2 x} \right| \\ \text{(1)} \quad &< \frac{\delta}{a^2 |x|} \end{aligned}$$

Our problem now is to obtain a constant upper bound for the factor $\frac{1}{a^2 |x|} = \frac{1}{a^2 |(x - a) + a|}$. It is sufficient to bound the denominator away from 0 or to guarantee

$$|x| = |(x - a) + a| > C > 0,$$

for some number C . We have (Appendix 1-3, Formula (3))

$$|x| = |(x - a) + a| \geq |a| - |x - a|.$$

Entering $|x - a| < \delta$ in this relation, we obtain

$$|x| \geq |a| - |x - a| > |a| - \delta.$$

To obtain a constant lower bound C we restrict* $\delta \leq \frac{|a|}{2}$. In that case

$$|x| > |a| - \delta > \frac{|a|}{2} > 0$$

and $C = \frac{|a|}{2}$.

It follows from $|x - a| < \delta$ that $|x| > \frac{|a|}{2}$ and $\frac{1}{|x|} < \frac{2}{|a|}$. (See Exercises A1-3, No. 20). Consequently, from (1), we have

$$\begin{aligned} |r(x) - L| &< \frac{\delta}{a^2|x|} \\ &< \delta \frac{2}{a^2|a|} \\ &\leq \delta \frac{2}{|a|^3} \end{aligned}$$

Step 2. The value of δ is restricted by two conditions:

$$\delta \leq \frac{|a|}{2} \quad \text{and} \quad \frac{2\delta}{|a|^3} \leq \epsilon.$$

To satisfy both conditions we take

$$\delta = \min\left\{\epsilon \frac{|a|^3}{2}, \frac{|a|}{2}\right\}.$$

*Of course, in general, we could have restricted δ in any convenient way so that $\delta < |a|$. For definiteness we took

$$\delta \leq \frac{|a|}{2}.$$

Step 3. Enter the above value of δ in the statement of the problem. The verification follows the pattern of Step 1. At the last line we use

$$\delta \leq \epsilon \frac{|a|^3}{2}$$

to obtain

$$|r(x) - L| < \epsilon.$$

Example 3-3f.

Statement of the Problem.

To prove for $a > 0$ that $\lim_{x \rightarrow a} r(x) = \lim_{x \rightarrow a} \frac{\sqrt{x} - \sqrt{a}}{x - a} = \frac{1}{2\sqrt{a}} = L.$

For each $\epsilon > 0$ obtain a δ .

Show: if $0 < |x - a| < \delta$, then $|r(x) - L| < \epsilon.$

(Observe that $r(x)$ is defined only for $x \geq 0$.)

Step 1.

$$\begin{aligned} (a) \quad r(x) - L &= \frac{\sqrt{x} - \sqrt{a}}{x - a} - \frac{1}{2\sqrt{a}} \\ &= \frac{1}{\sqrt{x} + \sqrt{a}} - \frac{1}{2\sqrt{a}} \quad (x \neq a) \\ &= \frac{\sqrt{a} - \sqrt{x}}{2\sqrt{a}(\sqrt{x} + \sqrt{a})} \\ &= \frac{a - x}{2\sqrt{a}(\sqrt{a} + \sqrt{x})^2} \end{aligned}$$

(Note that \sqrt{x} is not defined for negative values and therefore we guarantee $0 \leq x$ by imposing the restrictions $|x - a| < a$. For this purpose we require $\delta \leq a$.)

82

$$\begin{aligned}
 (b) \quad |r(x) - L| &= \left| \frac{a - x}{2\sqrt{a}(\sqrt{a} + \sqrt{x})^2} \right| \\
 &= \frac{\sqrt{|x - a|}}{2\sqrt{a}(\sqrt{a} + \sqrt{x})^2} \\
 &< \frac{\delta}{2\sqrt{a}(\sqrt{a} + \sqrt{x})^2}, \quad (\text{from } |x - a| < \delta) \\
 &\leq \frac{\delta}{2\sqrt{a}(\sqrt{a})^2} \quad (\text{from } \sqrt{x} \geq 0) \\
 &\leq \frac{\delta}{2(\sqrt{a})^3}
 \end{aligned}$$

Step 2. Take $\delta = \min\{2(\sqrt{a})^3\epsilon, a\}$.

Step 3. For the above value of δ every expression used in Step 1 is defined for all x in the deleted δ -neighborhood $0 < |x - a| < \delta$. (This requires $x \neq a$ and $x \geq 0$.) The verification follows Step 1. At the last line we use $\delta \leq 2(\sqrt{a})^3\epsilon$ to obtain

$$|r(x) - L| < \epsilon.$$

In the preceding examples we have not always followed the outline to the letter but used it only as a serviceable guide. Special difficulties are likely to appear in Step 1 and we cannot anticipate all contingencies. The only absolutely general pattern is the construction of a non-decreasing chain of expressions.

$$\phi_0 \leq \phi_1 \leq \phi_2 \leq \dots \leq \phi_n$$

where $\phi_0 = |r(x) - L|$, $\phi_n = g(\delta)$ and $\phi_1, \phi_2, \dots, \phi_{n-1}$ may involve both x and δ . To construct such a sequence in a particular case may require the greatest ingenuity.

In these examples we have verified that a given value L is actually the limit but have not shown how the limit L was obtained. In the next section we shall develop general theorems which will enable us to discover the value of the limit and to prove that the value is correct. Epsilonics will be necessary only to prove the theorems, not to apply them.

Exercises 3-3

1. Prove: $\lim_{x \rightarrow 4} (\frac{1}{2}x - 3) = -1$: obtain an upper bound $g(\delta)$ for the absolute error and find δ in terms of ϵ .

2. Give arguments that prove

(a) $\lim_{x \rightarrow a} c = c$, c any constant.

(b) $\lim_{x \rightarrow a} x = a$.

(c) $\lim_{x \rightarrow a} kx = ka$, k any constant.

(Use the results of Example 3-3a of the text for parts b and c.)

3. Invoke the definition directly to prove the existence of the limits in Problem 2.

4. In each of the following guess the limit, and then prove that your guess is correct.

(a) $\lim_{x \rightarrow 0} \frac{1}{1+x^2}$

(e) $\lim_{x \rightarrow 2} \frac{x^2 - 4}{x^3 - 8}$

(b) $\lim_{x \rightarrow 3} \frac{x^2(x-3)}{x}$

(f) $\lim_{x \rightarrow 0} \frac{x^3 - 3x - 1}{x + 2}$

(c) $\lim_{x \rightarrow a} \frac{x^3 - a^3}{x - a}$

(g) $\lim_{x \rightarrow 1} \frac{4x^2 - 3x - 1}{x + 2}$

(d) $\lim_{x \rightarrow 1} \frac{x+1}{x^2+1}$

3-4. Limit Theorems.

If the epsilonic definition of limit were required in every calculation with limits, the development of the calculus would be so disjointed and so overburdened with elaborate detail that it could only be mastered by a few devoted specialists. We need and we shall derive theorems that broadly cover most of the significant calculations with limits. In the end it will only be the exceptional cases for which epsilonic techniques are necessary.

The first general results apply to rational combinations of functions, that is, expressions formed from the functions of a given set by the rational operations of addition, subtraction, multiplication, and division. If each function of the given set has a limit as x approaches a , then the limit of any rational combination of these functions is the same rational combination of the corresponding limits (with divisions by zero excluded).

There are certain special rational combinations, called linear combinations, which recur often in different contexts. It is worth distinguishing them as a class because of their importance. A linear combination is built up by addition of functions and multiplication of functions by constants. Such a linear combination can be put in the form

$$\phi : x \rightarrow \phi(x) = c_1 f_1(x) + c_2 f_2(x) + \dots + c_n f_n(x)$$

where c_1, c_2, \dots, c_n are constants. In particular, a polynomial of degree less than or equal to n can be written in the form

$$\phi(x) = c_0 + c_1 x + c_2 x^2 + \dots + c_n x^n,$$

and may therefore be thought of as a linear combination of powers $1, x, x^2, \dots, x^n$.

The evaluation of the limit of a linear combination is an instructive instance of the general method of evaluating the limits of rational combinations:

Example 3-4a.

$$\begin{aligned} \lim_{x \rightarrow 4} (6\sqrt{x} + 5x + \pi) &= \lim_{x \rightarrow 4} 6\sqrt{x} + \lim_{x \rightarrow 4} 5x + \lim_{x \rightarrow 4} \pi \\ &= (\lim_{x \rightarrow 4} 6)(\lim_{x \rightarrow 4} \sqrt{x}) + (\lim_{x \rightarrow 4} 5)(\lim_{x \rightarrow 4} x) + \lim_{x \rightarrow 4} \pi \\ &= 6 \cdot \lim_{x \rightarrow 4} \sqrt{x} + 5 \cdot \lim_{x \rightarrow 4} x + \pi \end{aligned}$$

Note that in the example we have used three limit theorems without proof; in essence these are:

- (1) The limit of the sum of two functions is the sum of the limits.
- (2) The limit of the product of two functions is the product of the limits.
- (3) The limit of a constant is that constant.

Consider the statement

$$\lim_{x \rightarrow a} c = c$$

Note that the interpretations of c on the right and left of this equation are slightly different. On the left, c stands for $f(x)$, where

$$f : x \rightarrow c$$

and on the right c is the particular value assumed by the function for each value of x . With this in mind we have

THEOREM 3-4a. For a constant function $f : x \rightarrow c$,

$$\lim_{x \rightarrow a} f(x) = c$$

Proof. We have

$$|f(x) - c| = |c - c| = 0 < \epsilon,$$

for every positive ϵ and every choice of δ . (The constant function is a trivial case, of course, but we include it for completeness.)

THEOREM 3-4b. If $\lim_{x \rightarrow a} f(x) = L$, then for any constant c ,

$$\lim_{x \rightarrow a} c f(x) = c \lim_{x \rightarrow a} f(x) = cL.$$

Proof. We may assume $c \neq 0$, for if $c = 0$, the problem is reduced to that of Theorem 3-4a. Given any $\epsilon > 0$, we wish to make

$$|c f(x) - cL| < \epsilon$$

by restricting x to a deleted neighborhood

$$0 < |x - a| < \delta.$$

From the hypothesis we know that for any ϵ^* we can find a δ^* so that if

$$0 < |x - a| < \delta^*,$$

then

$$|f(x) - L| < \epsilon^*,$$

and

$$|c f(x) - cL| = |c| \cdot |f(x) - L| < |c| \epsilon^*.$$

Accordingly, we choose $\epsilon^* = \frac{\epsilon}{|c|}$, obtain the appropriate value δ^* for this ϵ^* , and set $\delta = \delta^*$.

In the following theorems we require that in some deleted neighborhood of a the domains of the functions entering the combination all coincide. This requirement eliminates nonsensical combinations such as $f(x) + g(x)$ when $f(x)$ is defined only for $x > a$ and $g(x)$ is defined only for $x < a$. The likelihood of ever making such a mistake is extremely small and therefore we do not mention this restriction on the functions explicitly in the statements or proofs of the theorems.

THEOREM 3-4c. If $\lim_{x \rightarrow a} f(x) = L$ and $\lim_{x \rightarrow a} g(x) = M$, then

$$\lim_{x \rightarrow a} [f(x) + g(x)] = L + M.$$

Proof. We must show that for any given $\epsilon > 0$ there is some δ such that

$$|f(x) + g(x) - (L + M)| < \epsilon$$

for all x in the common domain of f and g satisfying

$$0 < |x - a| < \delta.$$

From the hypothesis we know that for any positive ϵ_1 and ϵ_2 , no matter how small, we can find δ_1 and δ_2 such that

$$|f(x) - L| < \epsilon_1 \quad \text{when} \quad 0 < |x - a| < \delta_1,$$

$$|g(x) - M| < \epsilon_2 \quad \text{when} \quad 0 < |x - a| < \delta_2.$$

But

$$\begin{aligned} |f(x) + g(x) - (L + M)| &= |f(x) - L + g(x) - M| \\ &\leq |f(x) - L| + |g(x) - M|. \end{aligned}$$

To keep within the tolerance ϵ we can choose ϵ_1 and ϵ_2 to be any positive quantities whose sum is ϵ . For convenience, we fix

$$\epsilon_1 = \epsilon_2 = \frac{\epsilon}{2}.$$

Taking the appropriate values δ_1 , δ_2 for these values ϵ_1 , ϵ_2 , we set

$$\delta = \min\{\delta_1, \delta_2\}.$$

For this choice of δ , whenever

$$0 < |x - a| < \delta,$$

then

$$|f(x) + g(x) - (L + M)| < \frac{\epsilon}{2} + \frac{\epsilon}{2} \leq \epsilon.$$

Since a linear combination can be built up by successive operations of addition of two functions and multiplication by a constant, we obtain

Corollary. The limit of a linear combination of functions is the same linear combination of the limits of the functions; i.e., if

$$\lim_{x \rightarrow a} f_i(x) = L_i, \quad i = 1, 2, \dots, n$$

then

$$\begin{aligned} \lim_{x \rightarrow a} [c_1 f_1(x) + c_2 f_2(x) + \dots + c_n f_n(x)] &= c_1 \lim_{x \rightarrow a} f_1(x) + c_2 \lim_{x \rightarrow a} f_2(x) \\ &+ \dots + c_n \lim_{x \rightarrow a} f_n(x) = c_1 L_1 + c_2 L_2 + \dots + c_n L_n. \end{aligned}$$

The proof is left as an exercise.

For general rational combinations we have the further operations of multiplication and division.

Example 3-4b.

$$\begin{aligned} \lim_{x \rightarrow 4} \left[\frac{1}{x} - 2x^2 \sqrt{x} \right] &= \lim_{x \rightarrow 4} \frac{1}{x} - (\lim_{x \rightarrow 4} 2)(\lim_{x \rightarrow 4} x^2)(\lim_{x \rightarrow 4} \sqrt{x}) \\ &= \frac{1}{\lim_{x \rightarrow 4} x} - 2(\lim_{x \rightarrow 4} x)(\lim_{x \rightarrow 4} x)(\lim_{x \rightarrow 4} \sqrt{x}) \\ &= \frac{1}{4} - 2 \cdot 4 \cdot 4 \cdot 2 = -63 \frac{3}{4}. \end{aligned}$$

For $\phi(x) = \frac{1}{x} - 2x^2\sqrt{x}$ let us see in detail how ϕ can be built up in simple steps. We set

$$f_1(x) = \sqrt{x} ,$$

$$f_2(x) = xf_1(x) , \quad (\text{multiplication})$$

$$f_3(x) = xf_2(x) , \quad (\text{multiplication})$$

$$f_4(x) = -2f_3(x) , \quad (\text{multiplication})$$

$$f_5(x) = \frac{g_1(x)}{g_2(x)} , \quad (\text{division})$$

where

$$g_1(x) = 1$$

and

$$g_2(x) = x$$

and then,

$$\phi(x) = f_4(x) + f_5(x) , \quad (\text{addition})$$

It is, of course, tedious and unnecessary to decompose any rational combination into its elementary building blocks, but it is important to realize that it can be done and to know how to do it. (For example, it would be necessary to do so in writing computer programs.) In the process we have seen that to prove the general theorem concerning limits of rational combinations we now need to prove only the two special theorems for the limits of the product and quotient of two functions.

THEOREM 3-4d. If $\lim_{x \rightarrow a} f(x) = L$ and $\lim_{x \rightarrow a} g(x) = M$, then

$$\lim_{x \rightarrow a} [f(x) \cdot g(x)] = LM .$$

Proof. We wish to estimate the difference $f(x)g(x) - LM$, using the knowledge of the differences $f(x) - L$ and $g(x) - M$ given in the hypothesis. Now

$$\begin{aligned} f(x)g(x) - LM &= (f(x) - L)g(x) + L(g(x) - M) \\ &= (f(x) - L)(g(x) - M) + M(f(x) - L) + L(g(x) - M); \end{aligned}$$

hence,

$$(1) \quad |f(x)g(x) - LM| \leq |f(x) - L| \cdot |g(x) - M| + |M| \cdot |f(x) - L| + |L| \cdot |g(x) - M|.$$

From the hypothesis, we know that for any positive numbers ϵ_1 and ϵ_2 , there are corresponding controls δ_1 and δ_2 such that

$$|f(x) - L| < \epsilon_1 \quad \text{for } 0 < |x - a| < \delta_1,$$

$$|g(x) - M| < \epsilon_2 \quad \text{for } 0 < |x - a| < \delta_2.$$

Thus if we choose $\delta = \min\{\delta_1, \delta_2\}$, it will follow from (1) that when $0 < |x - a| < \delta$ then

$$(2) \quad |f(x)g(x) - LM| < \epsilon_1 \epsilon_2 + |M| \epsilon_1 + |L| \epsilon_2.$$

In order to keep from exceeding the tolerance ϵ we shall choose ϵ_1 and ϵ_2 so that

$$\epsilon_1 \epsilon_2 + |M| \epsilon_1 + |L| \epsilon_2 \leq \epsilon;$$

this will then determine our choice of δ_1 and δ_2 , and in turn that of δ . For convenience, we require that $\epsilon_1 = \epsilon_2 = v$ and that $v \leq 1$. Then

$$(3) \quad \epsilon_1 \epsilon_2 + |M| \epsilon_1 + |L| \epsilon_2 \leq v(1 + |L| + |M|).$$

We are now ready to choose v and verify (3). Let

$$(4) \quad v = \min \left\{ 1, \frac{\epsilon}{1 + |L| + |M|} \right\}.$$

Choose the corresponding δ_1 and δ_2 and let $\delta = \min\{\delta_1, \delta_2\}$. Then it follows from (2) and (4) when $0 < |x - a| < \delta$ that

$$|f(x)g(x) - LM| < v(1 + |L| + |M|) \leq \epsilon$$

as desired.

Since a polynomial $p(x)$ is a linear combination of powers, and powers are themselves products,

$$x^k = x \cdot x \cdots x \quad (k \text{ factors, } k \geq 1),$$

we can establish the following corollary.

Corollary. For any polynomial function p ,

$$\lim_{x \rightarrow a} p(x) = p(a).$$

The proof of this corollary is left as an exercise (Exercises 3-4, No. 2).

To prove the limit theorem for a quotient $\frac{f(x)}{g(x)}$, it is only necessary to prove the limit theorem for a reciprocal $\frac{1}{g(x)}$. The rule for general quotients then follows from

$$\frac{f(x)}{g(x)} = f(x) \left[\frac{1}{g(x)} \right].$$

First we prove a useful preliminary result.

Lemma 3-4. If $\lim_{x \rightarrow a} g(x) = M$ and $M > 0$, then there exists a neighborhood of a where $g(x) > 0$ for x in the domain of g .

Proof. Since g has the limit M at a , there is a δ -neighborhood of a wherein $g(x)$ is closer to M than to zero:

$$|g(x) - M| < \frac{M}{2}.$$

In this neighborhood,

$$\frac{3M}{2} > g(x) > \frac{M}{2} > 0.$$

If the function ϕ has a negative limit at $x = a$ then, upon applying Lemma 3-4 to the function $-\phi$, we see at once that $\phi(x)$ is negative in some deleted neighborhood of a . As further consequences of Lemma 3-4 we have the following two corollaries.

Corollary 1. If $\lim_{x \rightarrow a} g(x) = M$ and $M \neq 0$, then there exists a neighborhood of a where $|\frac{3M}{2}| > |g(x)| > |\frac{M}{2}|$ for x in the domain of g .

Corollary 2. A limit of a function whose values are nonnegative is nonnegative.

The proofs of these corollaries are left as exercises. (Exercises 3-4, No. 3)

THEOREM 3-4e. If $\lim_{x \rightarrow a} g(x) = M$ and $M \neq 0$, then

$$\lim_{x \rightarrow a} \frac{1}{g(x)} = \frac{1}{M}$$

Proof. We have

$$\begin{aligned} (2) \quad \left| \frac{1}{g(x)} - \frac{1}{M} \right| &= \left| \frac{M - g(x)}{Mg(x)} \right| \\ &= \frac{|g(x) - M|}{|M| \cdot |g(x)|} \end{aligned}$$

provided $g(x) \neq 0$. However, from Corollary 1 to Lemma 3-4 there is a δ -neighborhood of a wherein $|g(x)| > \frac{M}{2}$. Furthermore, for any ϵ^* the neighborhood can be taken so small that also

$$|g(x) - M| < \epsilon^*$$

From (2), therefore, we have

$$\begin{aligned} \left| \frac{1}{g(x)} - \frac{1}{M} \right| &= \frac{|g(x) - M|}{|M| \cdot |g(x)|} \\ &< \frac{\epsilon^*}{|M| \cdot \frac{|M|}{2}} \\ &\leq \frac{2\epsilon^*}{M^2} \\ &< \epsilon, \end{aligned}$$

where in the last line we have taken

$$\epsilon^* = \frac{M^2 \epsilon}{2}$$

To complete the proof we choose the value of δ appropriate to this ϵ^* .

Corollary 1. If $\lim_{x \rightarrow a} f(x) = L$ and $\lim_{x \rightarrow a} g(x) = M$ where $M \neq 0$, then

$$\lim_{x \rightarrow a} \frac{f(x)}{g(x)} = \frac{L}{M}$$

Corollary 2. If p and q are polynomials, and if $q(a) \neq 0$, then

$$\lim_{x \rightarrow a} \frac{p(x)}{q(x)} = \frac{p(a)}{q(a)}.$$

In connection with these corollaries, we observe that if $\lim_{x \rightarrow a} g(x) = 0$,

the quotient $\frac{f(x)}{g(x)}$ may still have a limit. Under these conditions,

$\lim_{x \rightarrow a} f(x) = 0$ is a necessary but not sufficient condition for existence of

$\lim_{x \rightarrow a} \frac{f(x)}{g(x)}$. The primary example is the derivative of a function expressed as

the limit of a ratio for which the numerator and denominator both approach zero. It is not possible to make any general statement about the existence of the limit for such cases; it is possible that $\lim_{x \rightarrow a} f(x) = 0$ and yet that

the limit of the quotient does not exist (for example, $\lim_{x \rightarrow 0} \frac{x}{x^2}$). (See Exercises 3-4, Nos. 14 and 15.)

In estimating $\lim_{x \rightarrow a} f(x)$ we can often bound f below and above by functions g and h which have limits as x approaches a . In that case we expect that the limit of f is bounded below and above by the limits of g and h . This result is a direct consequence of the following theorem.

THEOREM 3-4f. If $f(x) \leq g(x)$ in some deleted neighborhood of a , and $\lim_{x \rightarrow a} f(x) = L$ and $\lim_{x \rightarrow a} g(x) = M$, then $L \leq M$.

Proof. Since $g(x) - f(x)$ is nonnegative it follows that

$$\lim_{x \rightarrow a} [g(x) - f(x)] = M - L \geq 0.$$

(Theorem 3-4c and Corollary 2 to Lemma 3-4.)

Corollary 1. [Sandwich Theorem.] If

$$h(x) < f(x) < g(x)$$

in some deleted neighborhood of a , and if

$$\lim_{x \rightarrow a} h(x) = K \text{ and } \lim_{x \rightarrow a} g(x) = M,$$

93

then, if $\lim_{x \rightarrow a} f(x)$ exists,

$$K \leq \lim_{x \rightarrow a} f(x) \leq M.$$

Corollary 2 [Squeeze Theorem.] If $h(x) \leq f(x) \leq g(x)$ in some deleted neighborhood of a and if

$$\lim_{x \rightarrow a} h(x) = \lim_{x \rightarrow a} g(x) = M,$$

then

$$\lim_{x \rightarrow a} f(x) = M.$$

Exercises 3-4

1. Prove the corollary to Theorem 3-4c.

2. Prove the corollary to Theorem 3-4d.

3. Prove the corollaries to Lemma 3-4.

4. Prove the corollaries to Theorem 3-4e.

5. Find the following limits, giving at each step the theorem on limits which justifies it.

(a) $\lim_{x \rightarrow 3} (2 + x)$

(b) $\lim_{x \rightarrow -1} (5x - 2)$

(c) $\lim_{x \rightarrow 0} \left(\frac{a}{1 + |x|} - b\sqrt{|x|} \right)$, where a and b are constants.

(d) $\lim_{x \rightarrow a} (x^3 + ax^2 + a^2x + a^3)$, where a is constant.

6. Find the following limits, giving at each step the theorem which justifies it.

(a) $\lim_{x \rightarrow 1} \frac{x^3 - 1}{x^2 - 1}$

(b) $\lim_{x \rightarrow 3} \frac{x^2 - 9}{x^3 - 27}$

7. Find $\lim_{x \rightarrow 1} \frac{x^n - 1}{x - 1}$, for n a positive integer. Verify first that

$$\frac{x^n - 1}{x - 1} = x^{n-1} + x^{n-2} + \dots + x + 1, \quad (x \neq 1).$$

8. Determine whether the following limits exist and, if they do exist, find their values.

(a) $\lim_{x \rightarrow 1} \frac{1 + \sqrt{x}}{1 - x}$

(b) $\lim_{x \rightarrow a} (x^n - a^n)$; n is a positive integer, a is constant.

(c) $\lim_{x \rightarrow -1} \frac{\sqrt{2+x} + 1}{x+1}$

(d) $\lim_{x \rightarrow 1} \frac{(x-2)(\sqrt{x}-1)}{x^2 + x - 2}$

(e) $\lim_{x \rightarrow 1} \frac{1 - \sqrt{x}}{1 - x}$

9. Using the algebra of limits show that $\lim_{x \rightarrow a} \frac{f(x) - f(a)}{x - a} = L$ if and only

if $\lim_{x \rightarrow a} \frac{f(x) - f(a) - L(x - a)}{|x - a|} = 0$.

10. Assume $\lim_{x \rightarrow 0} \sin x = 0$ and $\lim_{x \rightarrow 0} \cos x = 1$. Find each of the following

limits, if the limit exists, giving at each step the theorem on limits which justifies it.

(a) $\lim_{x \rightarrow 0} \sin^3 x$

(d) $\lim_{x \rightarrow 0} \frac{\sin x}{\tan x}$

(b) $\lim_{x \rightarrow 0} \tan x$

(e) $\lim_{x \rightarrow 0} \frac{1 - \cos x}{\sin x}$

(c) $\lim_{x \rightarrow 0} \sin 2x$

(f) $\lim_{x \rightarrow 0} \frac{\cos 2x}{\cos x + \sin x}$

A11. (a) Prove Corollary 1 to Theorem 3-4f.

(b) Prove Corollary 2 to Theorem 3-4f.

(Hint: Prove $\lim_{x \rightarrow a} f(x)$ exists.)

A12. For what integral values of m and n does $\lim_{x \rightarrow a} \frac{x^m + a^m}{x^n + a^n}$ exist? Find the limit for these cases.

13. Prove that if $\lim_{x \rightarrow a} f(x) = 0$ and $g(x)$ is bounded in a neighborhood of

$x = a$, then $\lim_{x \rightarrow a} f(x) \cdot g(x) = 0$.

14. (a) Verify that if $\lim_{x \rightarrow a} \frac{f(x)}{g(x)}$ exists and if $\lim_{x \rightarrow a} g(x) = 0$, then

$\lim_{x \rightarrow a} f(x) = 0$.

(b) Describe functions f and g for which $\lim_{x \rightarrow a} f(x) = 0$ and

$\lim_{x \rightarrow a} g(x) = 0$ yet the limit of their quotient does not exist.

A15. Prove that if $\lim_{x \rightarrow a} g(x) = 0$ and $\lim_{x \rightarrow a} f(x)$ does not exist, then the

limit of the quotient $\frac{f(x)}{g(x)}$ does not exist.

16. The right-hand limit at a point $P(p, f(p))$ of a function is the limit of the function at the point P for a right-hand domain $(p, p + \delta)$. Similarly, for the left-hand limit, the domain is restricted to $(p - \delta, p)$. We denote them symbolically, by $\lim_{x \rightarrow p^+} f(x)$ and $\lim_{x \rightarrow p^-} f(x)$, respectively. In particular, $\lim_{x \rightarrow 2^+} [x] = 2$, $\lim_{x \rightarrow 2^-} [x] = 1$. Determine the indicated limits, if they exist, of the following:

(a) $\lim_{x \rightarrow 2^+} \frac{[x]^2 - 4}{x^2 - 4}$

(b) $\lim_{x \rightarrow 2^-} \frac{[x]^2 - 4}{x^2 - 4}$

(c) $\lim_{x \rightarrow 3^+} (x - 2 + [2 - x] - [x])$

(d) $\lim_{x \rightarrow 3^-} (x - 2 + [2 - x] - [x])$

(e) $\lim_{x \rightarrow 0^+} \left(\frac{x}{a} \left[\frac{b}{x} \right] - \frac{b}{x} \left[\frac{x}{a} \right] \right), a > 0, b > 0$

(f) $\lim_{x \rightarrow 0^-} \left(\frac{x}{a} \left[\frac{b}{x} \right] - \frac{b}{x} \left[\frac{x}{a} \right] \right), a > 0, b > 0$

(g) $\lim_{x \rightarrow 0^+} \frac{\sqrt{x}}{\sqrt{4 + \sqrt{x}} - 2}$

3-5. The Idea of Continuity.

The idea of continuity is intuitive: most of the functions which we studied in Chapter 2 had unbroken graphs, and these are appropriately called continuous. It is intuitive that a moving object traces a continuous path: the object cannot disappear at one place and reappear instantaneously at another. If we describe the path by $s = \phi(t)$, then the function ϕ is necessarily continuous.

We have also seen graphs with breaks or gaps, e.g., Figures 3-2b and 3-2c. Such graphs may be the appropriate way to represent certain kinds of physical situations. For instance, light moves through the air at a velocity of approximately $1.00c$, through water at a velocity of approximately $0.75c$, where c is the velocity of light in a vacuum. If we use a function $s \rightarrow v$ to describe the velocity at a distance a along a beam of light that penetrates a still body of water, then the function is discontinuous at the water surface.

The graphs of Section 3-2 show some ways in which a function can fail to be continuous, and they guide us to an informal definition of continuity. The function f_1 of Figure 3-2b is merely undefined at $x = 0$, and consequently its graph has a break there. The function f_2 has a value $f_2(0)$, but the point $(0, f_2(0))$ does not fill the gap in the graph: f_2 has the limit $L = 1$ as x approaches 0 , so that the values of f_2 near $x = 0$ are successively better approximations to L ; but $f_2(0)$ is not approximated by these functional values. For f_3 , on the other hand, the function is defined at $x = 0$ and the values of $f_3(x)$ do approximate $f_3(0)$ as x approaches 0 .

Previously, while investigating the limit of a function f as x approximates a , we paid no attention to the value of f at a , or even to the question whether f was defined at a . As abstract concepts, the value of f at a and the limit of f for x approximating a are unrelated. However, we need the concept of a continuous function; if $\lim_{x \rightarrow a} f(x) = L$ then the graph of f will have no gaps if $L = f(a)$.

We first met discontinuous functions in our differentiation procedure of Chapter 2. To find the derivative of $f : x \rightarrow x^2$, we need to find the limit as x approximates a of the ratio

$$r(x) = \frac{x^2 - a^2}{x - a} = \frac{x - a}{x - a} (x + a).$$

There is a gap in the domain of r : the formula for r is meaningless at $x = a$. To define the derivative we fill the gap. We observe that $r(x) = x + a$ when $x \neq a$ and that the value of $x + a$ when $x = a$ is the limit for $r(x)$. We use the limit to fill the gap in the values of r , and thus we replace the discontinuous function r by a continuous function ρ , where

$$\rho(x) = \begin{cases} r(x) & , x \neq a \\ \lim_{x \rightarrow a} r(x) & ; x = a \end{cases}$$

The function g_1 of Figure 3-2c illustrates the fact that it is not always possible to redefine a function so as to make it continuous at a point even when the point is contained in an interval of the domain of the function. To fill the gap in the domain of g_1 we would have to choose a value for $g_1(a)$ which is a limit at a for the function on both the restricted domains: one consisting of all x for which $x > a$ and the other of all x for which $x < a$. However, the two restrictions of the domain lead to different limits so that no single number is approximated as x approaches a . We are at liberty to redefine $g(a)$ as any real number we wish; but since g fails to have a limit for x approaching a , the values of g cannot serve as adequate approximations to $g(a)$ for all x in a deleted neighborhood of a , and therefore cannot go over continuously into the value $g(a)$.

From the study of these examples, we can abstract an informal definition of the concept: continuity of the function f at $x = a$. To ensure that the functional values have no break, f must satisfy three conditions:*

- (1) $f(a)$ exists;
- (2) $\lim_{x \rightarrow a} f(x)$ exists;
- (3) $f(a) = \lim_{x \rightarrow a} f(x)$.

Continuity will fail at $x = a$ if f is undefined (e.g., f_1 of Figure 3-2b at $x = 0$), if the limit fails to exist (e.g., g_1 of Figure 3-2c at $x = a$), or if functional value and limit both exist, but are not the same (e.g., f_2 of Figure 3-2b at 0). If the limit of f exists but the functional value is undefined, then a new function so defined as to agree with f at points other than $x = a$ and to have the value $\lim_{x \rightarrow a} f(x)$ at $x = a$ will be con-

*If condition (3) holds then conditions (1) and (2) are implied. We shall therefore adopt (3) as the basic definition of continuity.

tinuous at $x = a$ (e.g., f_1 of Figure 3-2b replaced by f_3). If the limit of $f(x)$ as x approaches a fails to exist, then any function agreeing with f in a deleted neighborhood of a is doomed to be discontinuous there (e.g., g_1 of Figure 3-2c at a).

Before establishing the properties of continuous functions, we supplement the preceding discussion with an analytical definition. For this purpose we express the relevant property, $\lim_{x \rightarrow a} f(x) = f(a)$, in epsilonic terms by using the Definition 3-2.

DEFINITION 3-5. The function f is defined to be continuous at a point a of its domain if every deleted neighborhood of a contains other points of the domain, and, for every $\epsilon > 0$, there exists a $\delta > 0$ such that

$$|f(x) - f(a)| < \epsilon$$

for every x in the domain of f that satisfies the inequality

$$|x - a| < \delta.$$

We verify that Definition 3-5 is equivalent to the earlier three-part definition. The value $f(a)$ is required to exist, since a must be in the domain of f . If we have $|f(x) - f(a)| < \epsilon$ for all x such that $|x - a| < \delta$, we surely have it for all x such that $0 < |x - a| < \delta$, so that $\lim_{x \rightarrow a} f(x)$ exists. Finally, because $f(a)$ appears in the inequality

$|f(x) - f(a)| < \epsilon$ of the present definition in precisely the place where $|f(x) - L| < \epsilon$ of Definition 3-2 we have $\lim_{x \rightarrow a} f(x) = f(a)$. All three ingredients of the three-part definition are implied by the formal one.

Conversely, we arrived at the formal definition by replacing the terms of the three-part definition with their analytical expressions; therefore, the validity of the informal definition implies that of the formal one. The two definitions for the continuity of a function at a point interior to an interval in its domain are equivalent.

A geometrical interpretation of Definition 3-5 is immediate: corresponding to any positive ϵ , there is an interval I of width 2ϵ with midpoint $f(a)$ such that for every $x \in I$, $f(x)$ is in an interval of width 2ϵ with $f(a)$ as midpoint. In other terms, for each ϵ it is possible to

confine the graph of $y = f(x)$ to the strip $f(a) - \epsilon < y < f(a) + \epsilon$ by restricting the domain of f to the interval I (Figure 3-5).

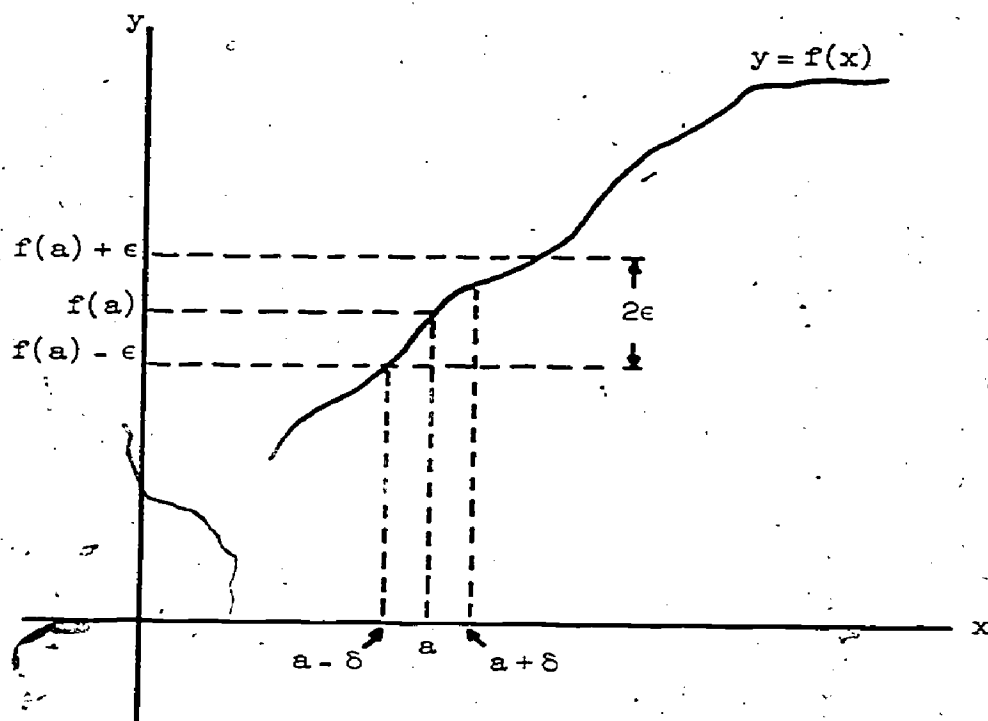


Figure 3-5

Note that to find a limit as x approaches a we examine deleted neighborhoods of the point, but to check continuity at a we include the point itself and examine an entire neighborhood, without deletion. It is because of this distinction that the inequality

$$0 < |x - a| < \delta$$

of Definition 3-2 may be replaced by the inequality

$$|x - a| < \delta$$

of Definition 3-5; moreover, the value $f(a)$ at $x = a$ now plays the role of the limit L .

If the condition of the definition is met, then a is a point of continuity for f . A point a where this condition is not met is a point of discontinuity of f , and f is discontinuous there.

Note particularly that if $f(x)$ is defined by a formula which has no real value when $x = a$, then f is discontinuous at a . For example, the function f given by $f(x) = \frac{1}{x-a}$ is discontinuous at a . (It is also true by this definition that the function g given by $g(x) = \sqrt{x^2 - 1}$ is discontinuous at each point a satisfying $|a| < 1$, but points separated from the domain of f are of no interest here.)

Example 3-5a. The absolute value function $f(x) = |x|$ (Figure 3-5), is defined by

$$f(x) = \begin{cases} x & \text{if } 0 \leq x, \\ -x & \text{if } x < 0. \end{cases}$$

We ask whether the function is continuous at $x = 0$, where the two intervals of definition meet. Given any $\epsilon > 0$, can we find $\delta > 0$ such that $|f(x) - f(0)| = ||x| - 0| < \epsilon$ for all x that satisfy the inequality $|x - 0| < \delta$? If $x > 0$, we wish to satisfy $|x| < \epsilon$ whenever $x < \delta$ while if $x < 0$ we wish to satisfy $|-x| < \epsilon$ whenever $-x < \delta$; in either case, the choice $\delta = \epsilon$ clearly is sufficient to hold $|x|$ within the tolerance ϵ , and establish the continuity.

Example 3-5b. From the graph of the integer part function $f : x \rightarrow [x]$ (Figure A2-1d) we expect to have discontinuity at the integers and continuity elsewhere. This is easily verified. If a is not an integer then $f(x) = [a]$ whenever $[a] \leq x < [a] + 1$. Let δ be the distance from a to the closer of the two integers $[a]$ or $[a] + 1$. If $|x - a| < \delta$ then $|[x] - [a]| = 0$ which is less than any error tolerance $\delta > 0$. On the other hand, for an integer n we have $f(n) = n$ and $f(x) = n - 1$ for $n - 1 \leq x < n$. Consequently, if $\epsilon < 1$, then no matter what neighborhood of n is taken there are always values of x on the left of n within the neighborhood for which

$$|[x] - [n]| = 1 > \epsilon$$

and the criterion of Definition 3-5 cannot be satisfied. (Note, however, that the function $g(x)$ defined on the restricted domain $n \leq x < n + 1$ by $g(x) = [x]$ is continuous at each point of its domain.)

Exercises 3-5

1. Use the formal definition of continuity to show that each of the following functions is continuous at $x = 1$.

$$(a) f : x \rightarrow \frac{x+1}{x^2+1}$$

$$(b) g : x \rightarrow \frac{4x^2 - 3x - 1}{x+2}$$

(See Exercises 3-3, Nos. 4d and 4g.)

2. For what value of x is each of the following functions discontinuous? Justify your answer.

$$(a) f : x \rightarrow \frac{x}{x}$$

$$(b) f : x \rightarrow \frac{x^2}{x+1}$$

3. Which of the following functions are discontinuous at $x = -1$? Justify your answer.

$$(a) f : x \rightarrow \frac{x+1}{x^3+1}$$

$$(b) g : x \rightarrow \frac{1}{1+x^2}$$

$$(c) h : x \rightarrow \frac{1}{1-x^2}$$

4. Discuss the points of discontinuity of $f : x \rightarrow [x] + [-x]$.
(See Exercises 3-1, No. 1.)

5. Prove that $f : x \rightarrow x - [x]$ is continuous for every x which is not an integer and discontinuous for integral values of x .

6. For each of the following functions define a new function which agrees with the given one for $x \neq a$ and is continuous at $x = a$.

$$(a) f : x \rightarrow \frac{x^3 - 1}{x^2 - 1}, a = 1$$

$$(b) f : x \rightarrow \frac{x^2 - 4}{x^3 - 8}, a = 2$$

$$(c) f : x \rightarrow \frac{x^n - 1}{x - 1}, a = 1, n \text{ an integer}$$

$$(d) f : x \rightarrow \frac{1 - \sqrt{x}}{1 - x}, a = 1$$

$$(e) f : x \rightarrow \frac{(x-2)(\sqrt{x}+1)}{x^2+x-2}, a = 1$$

7. For each of the following functions, if possible, define a new function which agrees with the given one for $x \neq 0$ and is continuous at $x = 0$. If this is not possible, state why.

(a) $f: x \rightarrow \frac{x}{x}$

(c) $f: x \rightarrow \frac{1}{|x|}$

(b) $f: x \rightarrow \frac{x}{|x|}$

(d) $f: x \rightarrow x - [x]$

8. For each of the following functions show that no function which agrees with the given one for $x \neq a$ can be so defined as to be continuous at $x = a$.

(a) $f: x \rightarrow \frac{1 + \sqrt{x}}{1 - x}$, $a = 1$

(b) $g: x \rightarrow \frac{\sqrt{2 + x + 1}}{x + 1}$, $a = -1$

9. For $x \neq 0$, let $f(x) = \left\lfloor \frac{1}{x} \right\rfloor$.

- (a) Sketch the graph of f over the closed intervals $[-2, -\frac{1}{6}]$ and $[\frac{1}{6}, 2]$.

- (b) Describe the behavior of $f(x)$ as x approximates 0 by positive values; by negative values.

- (c) Can you define a function which agrees with f for $x \neq 0$ and which is continuous at $x = 0$?

10. If f is an increasing function whose domain is the set of all real numbers, and if f is not continuous at a , what can you say of

$$\lim_{x \rightarrow a} |f(x) - f(a)| ?$$

11. For every real x , let $N(x)$ denote the number of distinct real square roots of x , i.e., the number of distinct real solutions of $y^2 = x$. Where does N have a limit? What is the limit? Where is N continuous? Let $P(x) = (N(x) - 1)^2$. Where does P have a limit? What is the limit? Where is P continuous? How does P differ from the function $f: x \rightarrow 1$? from the function $g: x \rightarrow \frac{x}{x}$?

12. Each of the functions f , g , and h is defined for all real x . Which of the functions is not continuous at 0 ?

$$f(x) = \begin{cases} 0, & x \text{ rational} \\ 1, & x \text{ irrational} \end{cases}$$

$$g(x) = \begin{cases} 0, & x \text{ rational} \\ x, & x \text{ irrational} \end{cases}$$

$$h(x) = \begin{cases} 1, & x \text{ rational} \\ -1, & x \text{ irrational} \end{cases}$$

13. Give an example of a function which is not continuous at 0 but whose absolute value is continuous at 0 .
14. (a) Show that the function f of No. 12 is periodic and determine all possible periods.
- (b) Show that every nonconstant periodic function which is continuous, at least at one point, has a fundamental (smallest) period.

15. If $f(x) = \begin{cases} \frac{1}{q}, & x = \frac{p}{q} \text{ rational; } p, q \text{ are relatively prime} \\ 0, & x \text{ irrational} \end{cases}$

show that f is continuous for all irrational x , and discontinuous for all rational x .

3-6. Properties of Functions Continuous at a Point.

(i) Rational combinations of continuous functions. We have proved that the elementary rational operations are preserved by the limit process: that is, that the limit of a sum is the sum of the limits of its terms; that the limit of a product is the product of the limits of its factors; and that the limit of a quotient is the quotient of the limits of its numerator and its denominator, provided that the limit of the denominator is not zero. It is immediate that if two functions are continuous at $x = a$ then so are their sum, product, and quotient if no division by zero is involved.

In the following theorems, as for the corresponding theorems on limits, we presuppose that the domains of the functions appearing in a combination coincide in some neighborhood of a .

THEOREM 3-6a. If the functions f and g are continuous at $x = a$, then so is the function h defined by $h(x) = f(x) + g(x)$. That is, the sum of two functions continuous at a point is also continuous there.

Proof. From Theorem 3-4c on the limit of a sum and from the definition of continuous functions we have

$$\begin{aligned}\lim_{x \rightarrow a} h(x) &= \lim_{x \rightarrow a} (f(x) + g(x)) \\ &= \lim_{x \rightarrow a} f(x) + \lim_{x \rightarrow a} g(x) \\ &= f(a) + g(a) \\ &= h(a).\end{aligned}$$

In precisely the same way, we obtain the following theorems.

THEOREM 3-6b. If the functions f and g are continuous at $x = a$, then so is the function h defined by $h(x) = f(x) \cdot g(x)$. That is, the product of two functions continuous at $x = a$ is also continuous there.

THEOREM 3-6c. If the functions f and g are continuous at $x = a$, and if $g(a) \neq 0$, then the function h , defined for $g(x) \neq 0$ by $h(x) = \frac{f(x)}{g(x)}$, is continuous at $x = a$. In other words, the quotient of continuous functions is continuous if no division by zero is involved.

From these results it can be proved that any rational combination of continuous functions is continuous at points where the denominator does not vanish. In particular, since constant functions are continuous (Theorem 3-4a), it can be proved that any linear combination of continuous functions is continuous. We have already seen that every polynomial is continuous (Corollary to Theorem 3-4a), hence, every rational function is continuous except when the denominator is zero.

Exercises 3-6a

1. Prove that $f : x \mapsto x^2$ is continuous at $x = a$, where a is any real number.
2. (a) Prove Theorem 3-6b using the limit theorems as in the proof of Theorem 3-6a.
(b) Prove Theorem 3-6c in the same way.
3. Prove Theorems 3-6b and 3-6c directly from Definition 3-5.
4. (a) If the function f is continuous at $x = a$ and the function g is not continuous at $x = a$, show that $f + g$ is not continuous at $x = a$.
(b) Can $f + g$ be continuous at $x = a$ if neither f nor g is continuous at $x = a$? Illustrate your answer by giving an example.
(c) Repeat the above using $f \cdot g$ for $f + g$.
5. Determine where the function $f : x \mapsto [x] + \sqrt{x - [x]}$ is continuous.

(ii) Continuity and differentiability. The functions which concern us in the calculus are usually continuous, but we shall generally not have to prove continuity as an independent fact. If a function is differentiable at a point, it is also continuous at the same point.

THEOREM 3-6d. If the function f has a derivative m at $x = a$, then f is continuous at $x = a$.

Proof. We are given that

$$\lim_{x \rightarrow a} \frac{f(x) - f(a)}{x - a} = m.$$

Furthermore,

$$\lim_{x \rightarrow a} (x - a) = 0.$$

Using Theorem 3-4d, we obtain

$$\begin{aligned} \lim_{x \rightarrow a} (f(x) - f(a)) &= \lim_{x \rightarrow a} \left[\frac{f(x) - f(a)}{x - a} \cdot (x - a) \right] \\ &= \left[\lim_{x \rightarrow a} \frac{f(x) - f(a)}{x - a} \right] \left[\lim_{x \rightarrow a} (x - a) \right] \\ &= m \cdot 0 \\ &= 0. \end{aligned}$$

Finally, from this result and Theorem 3-4c, we have

$$\lim_{x \rightarrow a} (f(x) - f(a)) = \lim_{x \rightarrow a} f(x) - \lim_{x \rightarrow a} f(a) = 0,$$

or
$$\lim_{x \rightarrow a} f(x) = f(a).$$

Therefore, $f(x)$ is continuous at a .

We have already proved (Example 3-3f) that \sqrt{x} has a derivative at a , for $a > 0$. As a consequence, we have

Corollary. The function $f : x \mapsto \sqrt{x}$ is continuous at $x = a$, for $a > 0$.

The relation between differentiability and continuity is a one-way affair: a function may be continuous at a point without being differentiable there.

It is easy to supply an example.

Example 3-6a. In Example 3-5a we demonstrated that $f: x \rightarrow |x|$ is continuous when $x = 0$. Does it have a derivative at $x = 0$? If we examine the graph (Figure A2-1c), it seems most unlikely, because such a derivative would assign a slope to the graph of f at the origin. The graph consists of two straight lines having the slopes -1 to the left and $+1$ to the right of the point, and it appears meaningless to talk of the direction of this graph at the point. The proof that $|x|$ does not have a derivative at $x = 0$ is left to you as an exercise (Exercises 3-6b, No. 1).

(iii) Composition of functions. The formation of rational combinations is one of the principal methods for the construction of more complicated mathematical functions from simpler ones, and we have proved that in general the existence of a limit and the property of continuity are preserved by rational combination. In the same spirit we now examine whether the existence of a limit and continuity are preserved by composition.

Example 3-6b. Let $f(x) = (1+x)^{1/2}$, $g(x) = \frac{1}{x}$, and $h(x) = \cos x$. We have many possible compositions of f , g , and h , such as

$$(1) \quad gf(x) = \frac{1}{(1+x)^{1/2}};$$

$$(2) \quad fgh(x) = \left(1 + \frac{1}{\cos x}\right)^{1/2};$$

$$(3) \quad gg(x) = \frac{1}{\left(\frac{1}{x}\right)}.$$

In composing functions, we must pay careful attention to their domains. Thus, in (1) of Example 3-6b, although $x = -1$ is perfectly satisfactory for the evaluation of $f(x)$, it is not permissible for the evaluation of $gf(x)$, since $f(-1) = 0$ is not in the domain of g . In a more complex composition such as (2), the difficulty may be somewhat disguised. Here we cannot choose any point x for which $-1 < \cos x < 0$; for example, $\frac{\pi}{2} < x < \pi$ or $\pi < x < \frac{3\pi}{2}$. For if x lies in such an interval then $\frac{1}{\cos x} < -1$, but the domain of f is restricted to numbers not less than -1 . Even such an

elementary composition as (3) may have hidden dangers: we are tempted to say that $gg(x) = x$, but strictly speaking the composition has the value x only for $x \neq 0$, because $g(0)$ is undefined and, therefore, $gg(0)$ is meaningless.

We shall consider functions which can be built up by successive composition. For such functions, we shall obtain a general theorem (under certain restrictions to be stated); namely, that any function constructed by successive compositions of continuous functions is itself continuous. It is sufficient to prove the theorem for a single composition.

THEOREM 3-6e. Let g be a function continuous at a whose range is contained in the domain of f . If f is continuous at $g(a)$, then the composite function $x \rightarrow fg(x)$ is continuous at a .

Proof. For any tolerance $\epsilon > 0$, we can choose a control $\sigma > 0$ such that

$$|f(u) - fg(a)| < \epsilon$$

for

$$|u - g(a)| < \sigma,$$

since f is continuous at $g(a)$. In particular, for $u = g(x)$, it follows that

$$(1) \quad |fg(x) - fg(a)| < \epsilon$$

whenever

$$(2) \quad |g(x) - g(a)| < \sigma.$$

We now regard σ as a prescribed tolerance for the approximation of $g(a)$ by $g(x)$, and because g is continuous at a , we can choose a control δ such that (2) holds for all x in the domain of g satisfying

$$|x - a| < \delta.$$

Since (1) holds for these values of x the composite function is continuous at a .

Example 3-6c. Given $f(u) = \frac{1}{u}$, under the valid assumption that $g(x) = \sin x$ is continuous for all x , we want to find all the discontinuities of fg (the cosecant function). We may eliminate from consideration all points where Theorem 3-6e is valid; that is, all points x where g is continuous and for which f is continuous at $g(x)$. The function g imposes no restriction on the domain of the composite function since it is continuous for all x . The function f is continuous except when $u = 0$; here f is undefined and consequently discontinuous. The values x for which $g(x) = 0$ will be points of discontinuity for fg because the composition is undefined there. We conclude that the cosecant function is discontinuous for $x = 0, \pm \pi, \pm 2\pi, \text{etc.}$ ---in brief, whenever x is an integral multiple of π .

(iv) Continuity of the inverse function. From the geometrical relation between the graph of a function and its inverse it may seem evident that continuity at a point of the domain of the function implies continuity at the corresponding point of the graph of its inverse; that is, if f is continuous at a then its inverse g is continuous at $b = f(a)$. This statement is not quite true* but the result can be assured by requiring an extra condition.

*Our geometrical intuition corresponds more closely to the concept of continuity on an interval (Section 3-7) rather than continuity at a point. If continuity at all points of an interval is violated, curious behavior is possible. For example, the function f given by

$$f(x) = \begin{cases} x, & \text{for } x \geq 0 \text{ and } x \text{ rational} \\ -x, & \text{for } x < 0 \text{ and } x \text{ irrational,} \end{cases}$$

is continuous and one-to-one. Nevertheless, its inverse is continuous nowhere but at $0 = f(0)$.

THEOREM 3-6f. Let f be an increasing (or decreasing) function and let g be its inverse. If f is continuous at a then g is continuous at $b = f(a)$.

Proof. Since f is increasing the existence of the inverse g is not at issue (see A2-4). To prove continuity of g at b we need to establish two things: (1) every deleted neighborhood of b contains points of the domain of g ; (2) given any ϵ -neighborhood of $a = g(b)$ the function g maps all the points of its domain within some δ -neighborhood of b into the ϵ -neighborhood of a .

Since f is continuous at a , for each neighborhood I of $b = f(a)$ there is a neighborhood J of a wherein f maps the points of its domain into I . Furthermore, any such J contains points of the domain of f other than a , and since the mapping is one-to-one it follows that I contains points of the range of f (which is the same as the domain of g) other than b .

We now have two cases to consider: either there are points of the domain of f within the ϵ -neighborhood of a on both sides of a or on only one side. If there are such points u, v on both sides of a , $u < a < v$, let δ be the distance from b to the closer of the points $f(u)$ or $f(v)$;

$$\delta = \min\{|f(u) - b|, |f(v) - b|\}.$$

Since the mapping g preserves order (or reverses order when f is decreasing) it follows that the part of the domain of g within the δ -neighborhood of b is mapped into the ϵ -neighborhood of a . Thus we have found a δ for each ϵ and continuity is proved. With slight modification the proof also applies to the one-sided case.

Exercises 3-6b

1. Prove that $f : x \rightarrow |x|$ does not have a derivative at $x = 0$.

2. Let $f : x \rightarrow x^n$, where n is a positive integer.

(a) Use the binomial theorem to expand $(x + h)^n$.

(b) From the result of (a) derive a formula for

$$r(x) = \frac{f(x+h) - f(x)}{h}$$

(c) From the result of (b) deduce that

$$m = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h} = nx^{n-1}$$

State which limit theorems you are using.

(d) Use Theorem 3-6d to show that f is continuous at $x = a$, where a is any real number.

3. Let $f(x) = \sqrt{x}$, $g(x) = \frac{1}{x-1}$, and $h(x) = \sin x$. Describe the domain of the function given by

(a) $fg(x)$.

(b) $gf(x)$.

(c) $hg(x)$.

(d) $gh(x)$.

(e) $hfg(x)$.

4. Assume that the functions $x \rightarrow \sin x$ and $x \rightarrow \cos x$ are continuous for all x . Find the discontinuities of the function given by

(a) $f(x) = \sin \frac{1}{x}$.

(b) $f(x) = \tan x$.

(c) $f(x) = \frac{1}{4 - 3 \sin^2 x}$.

(d) $f(x) = \sin \cos x$.

(e) $f(x) = \tan \frac{x^2 - 1}{x + 1}$.

(f) $f(x) = \tan \cos x - \cos \tan x$.

5. Prove, if $\lim_{x \rightarrow a} g(x) = b$ and f is continuous at b , then the composite function $x \rightarrow fg(x)$ has a limit as x approaches a and $\lim_{x \rightarrow a} fg(x) = f(\lim_{x \rightarrow a} g(x)) = f(b)$.

A 6.. Prove that if $\lim_{x \rightarrow a} f(x) = L$, and $\lim_{x \rightarrow a} g(x) = M$, then

$$\lim_{x \rightarrow a} \sqrt{(f(x))^2 + (g(x))^2} = \sqrt{L^2 + M^2}.$$

3-7. Properties of Functions Continuous on an Interval.

(i) Extreme and intermediate value theorems. We are interested in curves, in motions, in processes which are continuous, and we wish to represent them by functions. Therefore we are interested in functions which are continuous not merely at one point but at every point of an interval I . Such a function is said to be continuous on the interval I . Since the concept of continuity at a point given in Definition 3-5 is a local property of a function--a property which is determined by the values of the function within any neighborhood of the point, no matter how small--it is not at all obvious that the properties which intuition would ascribe to functions continuous on an interval can, in fact, be derived from the definition.

Many aspects of our intuitive picture of continuity are implicit in the precise definition. For example, we may think of the graph of a continuous function f passing through points $(a, f(a))$ and $(b, f(b))$ as the path of a walk over hilly terrain (Figure 3-7a).

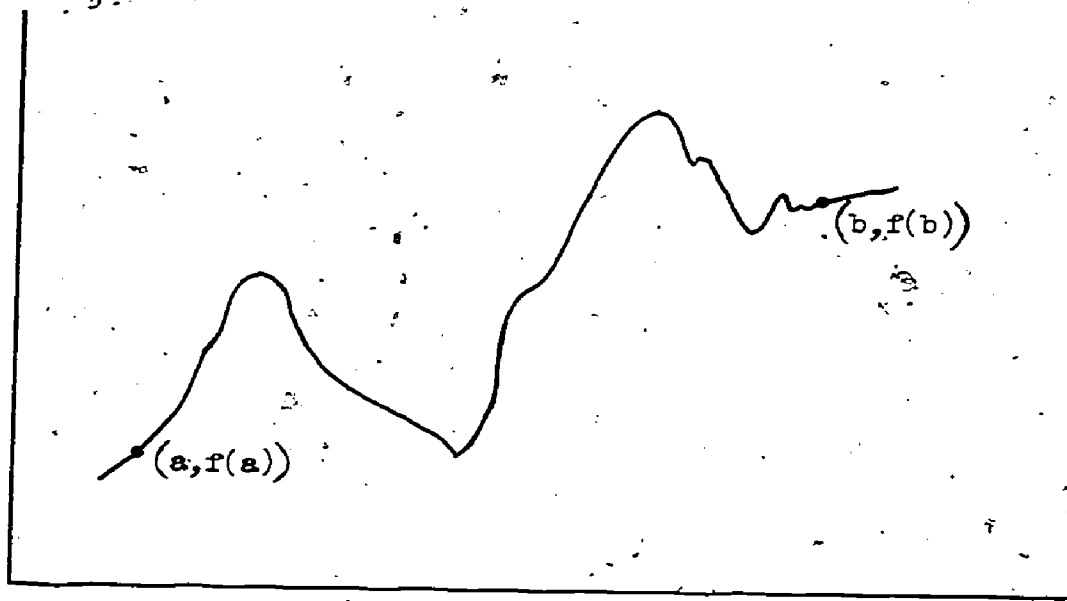


Figure 3-7a

It seems clear that such a path in passing from the elevation $f(a)$ to $f(b)$ must pass through every elevation between. Formally, this idea is expressed by the following theorem:

THEOREM 3-7a. (Intermediate Value Theorem). Let f be continuous on the closed interval $a \leq x \leq b$. Let v be any number between $f(a)$ and $f(b)$. Then there exists some value u in the interval such that $f(u) = v$.

Again, intuitively, there is at least one point on the path (possibly an end point) where the highest elevation on the interval is reached and another (or others) where the lowest elevation is reached. This property is expressed formally in the following theorem.

THEOREM 3-7b. (Extreme Value Theorem). If f is continuous on the closed interval $a \leq x \leq b$, then f has a maximum and a minimum in the interval. Specifically, within the interval, there exists at least one value x_M for which $f(x_M)$ is the maximum ($f(x_M) \geq f(x)$ for all x on the interval), and at least one value x_m for which $f(x_m)$ is the minimum ($f(x_m) \leq f(x)$ for all x on the interval).

The restriction to a closed interval is essential in Theorem 3-7b; e.g., the function given by $g(x) = x$ on the open interval $0 < x < 1$ has neither a maximum nor a minimum in the interval. The same remark applies to the function given by $h(x) = \frac{1}{x}$ on the same interval.

If $f(x_M)$ is the maximum of f on an interval, and $f(x_m)$ the minimum, then from the intermediate value theorem it is clear that f takes on every value between $f(x_M)$ and $f(x_m)$ on the interval between x_M and x_m . Hence we can combine the extreme and intermediate value theorems in a single statement:

Corollary. On a closed interval, the range of a continuous function contains a maximum, a minimum, and all values between.

This statement can be put more briefly as follows: A continuous function maps a closed interval onto a closed interval. The upper and lower endpoints of the image interval are, of course, the maximum and minimum of the function values.

It would require a substantial digression to prove these theorems in all generality; the proofs require a more lengthy exploration into the properties of real numbers than is appropriate here. For those interested, this material is relegated to Appendix 4.

For increasing or decreasing functions the property of continuity is equivalent to the property of the preceding corollary, namely,

THEOREM 3-7c. Let f be an increasing (or decreasing) function and let the domain of f be an interval I . If f is continuous on I then the range of f is an interval (by preceding corollary), and, conversely, if the range of f is an interval, then f is continuous.

Proof. For the proof of the converse we show that if f is monotone and the range R of f is in an interval then f is continuous. Let x_0 be an interior point of I . Since f is monotone and R is an interval, $f(x_0)$ is an interior point of R . Consequently, for any sufficiently small ϵ the values $f(x_0) - \epsilon$ and $f(x_0) + \epsilon$ are in R ; that is, there exist points c and d in I for which $f(c) = f(x_0) - \epsilon$ and $f(d) = f(x_0) + \epsilon$. Since f is monotone, if x is between c and d then $f(x)$ is between $f(x_0) - \epsilon$ and $f(x_0) + \epsilon$. Hence to assume

$$|f(x) - f(x_0)| < \epsilon$$

it suffices to require

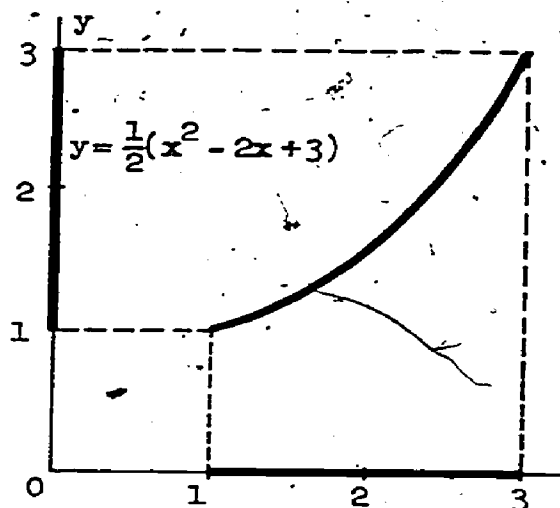
$$|x - x_0| < \delta,$$

where $\delta = \min\{|x_0 - c|, |x_0 - d|\}$.

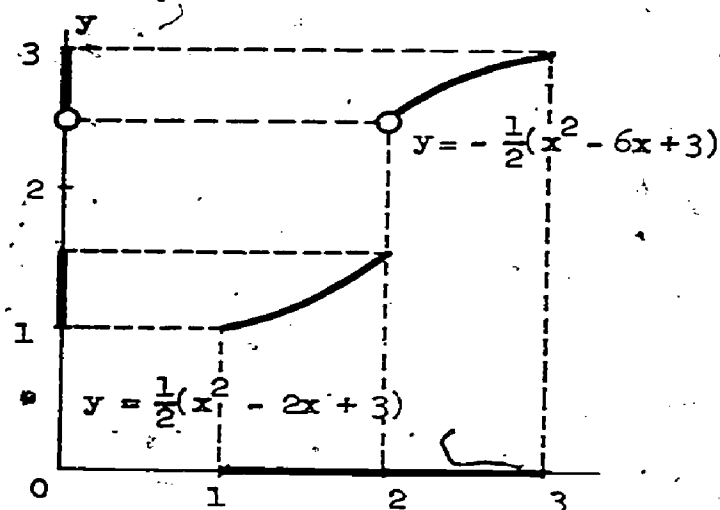
A slight modification of this argument suffices to prove the result where x_0 is an endpoint of I . Here we give an illustration. In Figure 3-7b(i) we depict the graph of a continuous monotone function $f: x \rightarrow \frac{1}{2}(x^2 - 2x + 3)$ for $1 \leq x \leq 3$; in (ii), the graph of a discontinuous increasing function g given by

$$g(x) = \begin{cases} \frac{1}{2}(x^2 - 2x + 3), & 1 \leq x \leq 2 \\ -\frac{1}{2}(x^2 - 6x + 3), & 2 < x \leq 3 \end{cases}$$

In both cases the projection of the graph on the x -axis, the domain of the function, is the interval $1 \leq x \leq 3$. For the continuous function f , the



(i)



(ii)

Figure 3-7b

projection on the y -axis, the range of the function, is an interval, $1 \leq y \leq 3$. For the discontinuous function g , the range consists of two separated intervals, $1 \leq y \leq \frac{3}{2}$ and $\frac{5}{2} < y \leq 3$. This situation is typical.

In the preceding theorems we have evidence that the formal definition of continuity on an interval does agree with various aspects of the intuitive conception of continuity. Yet, it is important to know that a precise, formal definition does not necessarily correspond in every respect to the intuitive idea from which it springs. The idea of continuity is a particularly revealing example. You have seen that a continuous function need not have a derivative at every point of its domain (e.g., $x \rightarrow |x|$). It is not obvious, but it is true, that there are functions continuous on an interval which do not have a derivative at any point of their domains. If you think of functions in terms of graphs which can be plotted, then this fact may be surprising; certainly no pen can follow the infinitely sinuous wiggles of such a graph. The presentation of such functions by the 19th century mathematician Weierstrass was, in fact, a definite shock to the mathematical world of his time. (An example of such a function is given in Appendix A4-3.)

In the early 19th century mathematicians evidently believed that a continuous function must also have a derivative, except perhaps at isolated points. In part this feeling probably stemmed from the earlier concept of function as a relation defined by a formal expression; the idea of function as we know it (then called "single-valued" function) had not been thoroughly explored. The formal expressions familiar to the mathematicians of that time would not be likely to suggest the peculiarities of nondifferentiable continuous functions. In 1872 Weierstrass proved definitely that

a certain function was just such a nondifferentiable continuous function. There is a persistent misconception that the mathematical community was deeply shocked by the example of Weierstrass. In fact, in 1834, thirty-six years prior to the paper of

Weierstrass, Bolzano^{1,2} gave an example of such a function. However, Bolzano did not completely prove that his function had all of the required properties. The story of the rediscovery of Bolzano's example and a proof is given by G. Kowalewski.² According to Weierstrass, around 1861 Riemann also proposed an example of such a function.³ Weierstrass found it too difficult to verify that the example of Riemann is correct, and it is unclear whether Riemann was able to verify this. Thus Weierstrass was the first to prove that a specific example actually was nowhere differentiable.

In Section A4-3 we present an example of a continuous but nowhere differentiable function whose properties can be demonstrated by the most elementary means.

-
1. C. B. Boyer, Concepts of the Calculus.
 2. G. Kowalewski, "Über Bolzano's Nichtdifferenzierbare Stetige Funktion." Acta Mathematica, 44. 1923, pp. 315-319.
 3. Weierstrass, Werke, Vol. 2, pp. 71-74.

Exercises 3-7

1. Exhibit a discontinuous function for which the range of the function is an interval.
2. On which of the following intervals is the function $f : x \rightarrow \sin x$ increasing? decreasing? In each case locate the maximum and minimum values if any.
 - (a) $(-\frac{\pi}{2}, \pi)$
 - (b) $[-\frac{\pi}{2}, \pi]$
 - (c) $[-\frac{\pi}{2}, 0]$
 - (d) $(-\frac{\pi}{2}, \frac{\pi}{2})$
 - (e) $[-\pi, \pi]$
3. Prove that $f : x \rightarrow x^n$ is increasing for $x > 0$ and ranges over the positive reals (n rational).
4. (a) Prove that $x \rightarrow \sqrt{f(x)}$ is continuous and increasing wherever f is positive, continuous, and increasing.
 (b) Prove that $x \rightarrow \sqrt[3]{f(x)}$ is continuous wherever f is continuous.
5. Prove that $\lim_{x \rightarrow 0} \frac{\sqrt{1+x} - \sqrt{1-x}}{x} = 1$.
6. Assume that $f : x \rightarrow \sin x$ is continuous for all x in $[-\frac{\pi}{2}, \frac{\pi}{2}]$.
 (a) Prove that f is increasing. Hint: Use the identity

$$\sin y - \sin x = 2 \sin \left(\frac{y-x}{2} \right) \cos \left(\frac{y+x}{2} \right)$$
 and consider as separate cases the closed intervals $[-\frac{\pi}{2}, 0]$ and $[0, \frac{\pi}{2}]$.
 (b) Show that the range of f is an interval. Show that the domain of the inverse of $f : x \rightarrow \sin x$ is an interval.

7. Assume that the function $f: x \rightarrow \tan x$ is continuous on the closed interval $[-\frac{\pi}{4}, \frac{\pi}{4}]$ and show that there is some number x where $-\frac{\pi}{4} < x < \frac{\pi}{4}$ such that $\tan x = \frac{\pi}{4}$. Can there be more than one such number?
8. Prove that if f is continuous on the closed interval $[a, b]$ and all values of f are in $[a, b]$, then there is an x in $[a, b]$ for which $f(x) = x$.
9. In which of the following intervals does $f: x \rightarrow |x|$ have a maximum? a minimum? Justify your answer.
- (a) $-1 < x < 1$
 - (b) $-1 \leq x < 1$
 - (c) $-1 < x \leq 1$
 - (d) $-1 \leq x \leq 1$
10. Show that the equation $x^4 + x - 10 = 0$ has at least one solution between $x = -2$ and $x = -1$ and obtain an approximation to the solution within a tolerance of $\frac{1}{2}$.
11. Isolate each real root of the given equation by exhibiting an interval containing this root and no others. (Each equation has four roots.)
- (a) $x^4 - x^3 - 9x^2 + 2x + 14 = 0$
 - (b) $2x^4 + 2x^3 - 3x^2 - x + 1 = 0$
12. (a) Show that the equation $\cos^2 x = \sqrt{|x|}$ has at least one positive root x , where $x < \frac{\pi}{6}$.
- (b) Find the maximum and the minimum value of $f: x \rightarrow \cos^2 x = \sqrt{|x|}$ on the closed interval $[0, \frac{\pi}{2}]$.
13. Prove that if $p(x)$ is a polynomial of odd degree (with real coefficients), then the equation $p(x) = 0$ has a real root.
14. Prove that the equation $x^n = a$ has exactly one negative root if n is an odd positive integer and $a < 0$.
15. Sketch the curves $y = x^n$ and $y = x^{1/n}$, $n = 0, 1, 2, \dots$ using the same set of axes.

16. Prove that if f is continuous and has an inverse on an interval then f is strongly monotone on the interval.
17. The temperature at any point of a thin circular ring is a continuous function of the point's position. Show that there is a pair of antipodes (points at opposite ends of a diameter) having the same temperature.
18. Sketch the graphs and determine how many points of discontinuity there are in the interval $[0, 2\pi]$ of the following functions:
- (a) $f : x \mapsto [\sin x]$.
 - (b) $f : x \mapsto [2 \sin x]$.
 - (c) $f : x \mapsto [a \sin x]$.
- A19. If f is periodic with periods 1 and $\sqrt{2}$ (i.e., $f(x) = f(x+1)$ and $f(x) = f(x+\sqrt{2})$) for all x , and if there is at least one point of continuity of f , show that f must be constant.
- A20. If g is continuous with $g(0) = g(1) = 1$ and, in the interval $[0, 1]$, $g(x^2) = (g(x))^2$, show that $g(x) = 1$ in $[0, 1]$.
21. The real roots of the equation $x^n + ax + b = 0$ (n a positive integer) can be "determined" by finding the intersections of the curves

$$y = x^n$$

and

$$y = -ax - b.$$

Verify the following table for the number of real roots of $x^n + ax + b = 0$.

- | | |
|-------------------------|--|
| (a) If n is even, and | $\begin{cases} b > 0, & \text{there are two or none,} \\ b < 0, & \text{there are two.} \end{cases}$ |
| (b) If n is odd, and | |
| | $\begin{cases} a > 0, & \text{there is one,} \\ a < 0, & \text{there are three or one.} \end{cases}$ |

Give numerical examples to illustrate each of the four cases.

Chapter 4

DIFFERENTIATION

4-1. Introduction.

The value of the concept of derivative is greatly enhanced by the simplicity of the techniques for computing the derivatives of common functions and functions constructed as combinations of them. In this chapter we shall find the derivatives of powers and trigonometric functions and show how to obtain the derivatives of functions obtained by rational combinations, inversions, and composition of functions with known derivatives.

By the definition in Chapter 2 the derivative of a function f at a point x is the limit as z approaches x of the ratio

$$r(z) = \frac{f(z) - f(x)}{z - x}.$$

From this we obtain the function

$$f' : x \longrightarrow \lim_{z \rightarrow x} \frac{f(z) - f(x)}{z - x},$$

which has as its domain the subset of the domain of f for which the limit exists. The function f' is called the derivative of f , that is, the function derived from f . Thus $f'(x)$ is the slope of the graph of f at the point x .

It is convenient for the purposes of computation to replace the denominator in the expression for $r(z)$ with a single letter h . We replace z by $x + h$ and obtain

$$f'(x) = \lim_{h \rightarrow 0} \frac{f(x + h) - f(x)}{h}.$$

When we are concerned with specific functions like x^2 or $\sin x$, the notation $f'(x)$ becomes somewhat awkward and we shall find it convenient to use the prefix D_x . Thus, for

$$f : x \longrightarrow x^2 \quad \text{and} \quad g : x \longrightarrow \sin x$$

we have

$$f'(x) = D_x x^2 \quad \text{and} \quad g'(x) = D_x \sin x.$$

The subscript is usually omitted when its reference is clear. We shall usually write Dx^2 and $D \sin x$, but when there may be doubt, we shall always use the subscript, e.g., $D_x \sqrt{x^2 + z^2}$.

Example 4-1a. For $f : x \rightarrow x^2$ we have

$$\begin{aligned} \frac{f(x+h) - f(x)}{h} &= \frac{(x+h)^2 - x^2}{h} \\ &= \frac{h}{h}(2x+h) \end{aligned}$$

Consequently, by the Corollary to Theorem 3-4c, we have

$$f'(x) = Dx^2 = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h} = 2x.$$

We list for ready reference the few derivatives which are given or follow easily from examples here and in Chapters 2 and 3. It will be the beginning of a catalog of functions whose derivatives you should learn. We shall add to the catalog in the following sections.

$$(1) \quad D_x C = 0, \quad C \text{ a constant.}$$

$$(4) \quad D\sqrt{x} = \frac{1}{2\sqrt{x}}.$$

$$(2) \quad Dx = 1.$$

$$(5) \quad D\left(\frac{1}{x}\right) = -\frac{1}{x^2}.$$

$$(3) \quad Dx^2 = 2x.$$

$$(6) \quad D|x| = \begin{cases} 1, & x > 0, \\ -1, & x < 0. \end{cases}$$

Exercises 4-1

1. Using the definition of the derivative, find $f'(x)$, where $f(x)$ equals:

(a) $2x^2 - x + 4$.

(b) $1 - x^3$.

(c) $\frac{1}{\sqrt{x}}$.

(d) $\frac{1}{x^2}$.

(e) $x^3 - 3x + 4$.

(f) $\frac{1}{ax + b}$.

(g) $ax^2 + bx + c$.

(h) $|x - 1|$.

(i) $ax + \frac{b}{x}$.

2. Given the line $y = 3x + 2$, find its slope at the points $(0, 2)$, $(-2, -4)$, and $(2, 8)$.

3. If $f(x) = 1 + 2x - x^2$, find the slope of the graph of f at points corresponding to:

(a) $x = 0$.

(c) $x = 1$.

(b) $x = \frac{1}{2}$.

(d) $x = -10$.

4. If $f(x) = x^3 + 2x + 1$, find all x such that

(a) $f'(x) = 0$.

(c) $f'(x) = 4$.

(b) $f'(x) = -1$.

(d) $f'(x) = 20$.

4-2. Rational Combinations.

We can extend the list of known derivatives greatly by considering rational combinations of differentiable functions. We need consider only three basic kinds of combinations, namely linear combinations, products, and quotients.

Linearity of differentiation.

THEOREM 4-2a. (Linearity of differentiation.) If f and g are both differentiable at the point x , then for any constants a and b we have

$$D_x (af(x) + bg(x)) = a Df(x) + b Dg(x).$$

Proof. By linearity of the operation of taking a limit (Corollary to Theorem-3-4c) we have

$$\begin{aligned} & \lim_{h \rightarrow 0} \left(\frac{af(x+h) - af(x)}{h} + \frac{bg(x+h) - bg(x)}{h} \right) \\ &= a \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h} + b \lim_{h \rightarrow 0} \frac{g(x+h) - g(x)}{h} \\ &= a Df(x) + b Dg(x). \end{aligned}$$

Example 4-2a. $D\left(3x^2 + 2\sqrt{x} - \frac{1}{x}\right) = 3D x^2 + 2D\sqrt{x} - D\left(\frac{1}{x}\right) = 6x + \frac{1}{\sqrt{x}} + \frac{1}{x^2}.$

Exercises 4-2a

1. Evaluate

(a) $D(4|x| + 6\sqrt{x})$.

(b) $D(5x^2 + \frac{2}{x})$.

(c) $D(|7x| + \frac{3}{2x-1})$, $x < \frac{1}{2}$.

(d) $D(|ax| - |bx|)$, $a > 0$, $b < 0$.

2. Consider $g: x \rightarrow |x+2| - |3-x|$.

(a) Sketch the graph of g .

(b) Define $g(x)$ explicitly in terms of linear functions for all real x .

(c) For what values of x is the derivative not defined?

3. Consider $f : x \rightarrow [x]$ ($[x]$ is the integer part of x , defined in A2-1.)
- (a) Find $f'(x)$ if it exists, at each of the values $x = -2.8$, $x = 0.6$, $x = 2$.
- (b) Find the domain of the derivative f' .
4. Consider $f : x \rightarrow x - [x]$.
- (a) Draw the graph of f .
- (b) Find $f'(-1.5)$ and $f'(2.3)$ and describe the domain of the derivative.
5. Extend Theorem 4-2a to a general linear combination of functions

$$\phi : x \rightarrow c_1 f_1(x) + c_2 f_2(x) + \dots + c_n f_n(x).$$

6. For each of the following functions, find the derivative and describe the domain of the derivative.
- (a) $f : x \rightarrow |x^2 - 2|$ (e) $f : x \rightarrow \operatorname{sgn}(1 - \sqrt[3]{x})$
- (b) $f : x \rightarrow [2x^2]$ (f) $f : x \rightarrow \max\{x^3, 4|x|\}$
- (c) $f : x \rightarrow |[x + 1]|$ (g) $f : x \rightarrow \min\{[x], \max\{x^3, 2x^2\}\}$
- (d) $f : x \rightarrow [|x + 1|]$ (h) $f : x \rightarrow \operatorname{sgn}(\min\{x^3 - 1, 7\})$

7. Right-hand and left-hand derivatives are defined in terms of right-hand and left-hand limits (see Exercises 3-4, No. 16) as follows:

Right-hand derivative: $D^+f(x) = \lim_{h \rightarrow 0^+} \frac{f(x+h) - f(x)}{h}$

Left-hand derivative: $D^-f(x) = \lim_{h \rightarrow 0^-} \frac{f(x+h) - f(x)}{h}$

In particular, $D^+|x| = D^-|x| = 1$, $x > 0$,

$D^+|x| = D^-|x| = -1$, $x < 0$,

$D^+|x| = 1 = -D^-|x|$, $x = 0$.

- (a) Show that $D^-|x^5| = D^+|x^5|$ for all x .
- (b) For what values of x does $D^-|x^3 - 2| = D^+|x^3 - 2|$?
- (c) Show that a function is differentiable at a point if and only if it has equal right-hand and left-hand derivatives at the point.

Derivatives of products.

THEOREM 4-2b. If the functions f and g are differentiable at x , then the product function

$$F : x \longrightarrow f(x)g(x)$$

has the derivative at x given by

$$F'(x) = f(x)g'(x) + g(x)f'(x).$$

Proof. From the definition of F we have

$$\begin{aligned} \frac{F(x+h) - F(x)}{h} &= \frac{f(x+h)g(x+h) - f(x)g(x)}{h} \\ &= \frac{f(x+h)g(x+h) - f(x+h)g(x) + f(x+h)g(x) - f(x)g(x)}{h} \\ &= f(x+h) \frac{g(x+h) - g(x)}{h} + g(x) \frac{f(x+h) - f(x)}{h}. \end{aligned}$$

Since f is differentiable at x , it is continuous there (Theorem 3-6d) and we have

$$\lim_{h \rightarrow 0} f(x+h) = f(x),$$

$$\lim_{h \rightarrow 0} g(x) = g(x),$$

$$\lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h} = f'(x),$$

$$\lim_{h \rightarrow 0} \frac{g(x+h) - g(x)}{h} = g'(x).$$

It follows from the theorems on limits (Section 3-4) that

$$\lim_{h \rightarrow 0} \frac{F(x+h) - F(x)}{h} = F'(x) = f(x)g'(x) + g(x)f'(x).$$

Example 4-2b. $D_x x^{3/2} = D(x\sqrt{x}) = xD\sqrt{x} + \sqrt{x} D_x x$

$$= x \frac{1}{2\sqrt{x}} + \sqrt{x} \cdot 1$$

$$= \frac{\sqrt{x}}{2} + \sqrt{x} = \frac{3}{2}\sqrt{x}$$

Corollary. If f' exists and if $F(x) = [f(x)]^2$, then

$$\begin{aligned} F'(x) &= 2[f(x)] D_x f(x) \\ &= 2f(x) \cdot f'(x). \end{aligned}$$

The proof is left as an exercise (Exercises 4-2b, No. 3).

Exercises 4-2b

1. Find the derivatives of the following functions.

(a) $x(2x - 3)$

(b) $(4x - 2)(4 - 2x)$

(c) $(x^2 + x + 1)(x^2 - x + 1)$

(d) $\sqrt{x} (ax + b)^3$

(e) $\frac{1}{x} \cdot \sqrt{x}$

(f) $\frac{1}{x} \cdot (5x + 2)$

(g) $|x| \cdot \frac{1}{x^2}$

(h) $x^{7/2}, x \neq 0$

(i) $3x^4 - \frac{1}{\sqrt{x}}$

(j) $3x^2(x^2 - 5)$

(k) $|13x^2 - 36 - x^4|$

(l) $|5x^2 - 36 + x^4|$

2. Evaluate

(a) $D(3x^2 + 5x - 1)^2$

(b) $D(3 - 5x)^3$

(c) $D(3 - 5x)^4$

(d) $D(x(\sqrt{x} - 1)^2)$

(e) $D(x + \frac{1}{x})^2$

(f) $D(\frac{x^{3/2}}{3} - \frac{x^{1/2}}{2} + x^{-1/2})$ (Hint: $x^{-1/2} = \frac{1}{x} \cdot \sqrt{x}$)

(g) $D(4\sqrt{x^3} - 2\sqrt{x} + \frac{1}{\sqrt{x}})$

(h) $D([x^2 - 1](x^2 - 3x + 1))$

(i) $D(x([x + 1] - 1))^n$, n an integer.

3. Prove the Corollary to Theorem 4-2b.

4. Find the listed derivative by two methods; first, expand and then differentiate; second, use the product formula.

(a) $D(x^2 + 1)^2$

(b) $D(x^2(x^2 + 1)^2)$

(c) $D((x + 1)(x^2 - x + 1))$

(d) $D((ax^2 + bx + c)(dx^2 + ex + f))$

5. Find the derivatives of each of the following functions in as many ways as you can and describe its domain. (Do not overlook the definition of the derivative.)

(a) $f : x \rightarrow x|x|$

(f) $f : x \rightarrow |x|[x]^2$

(b) $f : x \rightarrow |x|^2$

(g) $f : x \rightarrow x[4 - x^2]$

(c) $f : x \rightarrow x[x]$

(h) $f : x \rightarrow x \max\{x, 2 - x^2\}$

(d) $f : x \rightarrow [x]^2$

(i) $f : x \rightarrow x[x]|x|$

(e) $f : x \rightarrow |x|[x]$

(j) $f : x \rightarrow [|x|] \cdot |[x]|$

The corollary to Theorem 4-28 may be extended to any higher integral power of a function f as indicated in the exercises, Exercises 4-2c. We state it here for completeness.

THEOREM 4-2c. [Power Rule for Positive Integer n] If f' exists and if $F(x) = [f(x)]^n$, then

$$F'(x) = n[f(x)]^{n-1} f'(x)$$

for any positive integer n .

Example 4-2c. If $F(x) = (3x - 2)^5$, then

$$\begin{aligned} F'(x) &= 5(3x - 2)^4(3) \\ &= 15(3x - 2)^4. \end{aligned}$$

Corollary 1. If $G(x) = x^n$, then $G'(x) = nx^{n-1}$ for any positive integer n .

The proof is left as an exercise.

Since a polynomial

$$p(x) = a_0 + a_1x + a_2x^2 + \dots + a_nx^n$$

is a linear combination of powers of x , Theorem 4-2a and Corollary 1 enable us to differentiate any polynomial.

Corollary 2. A polynomial function p , where

$$p(x) = a_0 + a_1x + a_2x^2 + \dots + a_nx^n,$$

has a derivative for each real x given by

$$p'(x) = a_1 + 2a_2x + \dots + na_nx^{n-1}.$$

Example 4-2d. $D(4x^3 - 7x^2 + 3x - 2) = 12x^2 - 14x + 3$.

Theorems 4-2a and 4-2c enable us to differentiate polynomial functions of any function whose derivative we know.

Corollary 3. If p is a polynomial and if f' exists, then

$$D_x p(f(x)) = p'(f(x)) f'(x)$$

Example 4-2e.

$$\begin{aligned}
 (1) \quad D[x^{7/2} + 4x^{5/2}] &= D[(\sqrt{x})^7 + 4(\sqrt{x})^5] \\
 &= 7(\sqrt{x})^6 D\sqrt{x} + 20(\sqrt{x})^4 D\sqrt{x} \\
 &= 7x^3 \frac{1}{2\sqrt{x}} + 20x^2 \frac{1}{2\sqrt{x}} \\
 &= \frac{7}{2}x^{5/2} + 10x^{3/2}
 \end{aligned}$$

$$\begin{aligned}
 (2) \quad D_x[(3x-2)^5 + 3(3x-2)^4 - 5(3x-2)] \\
 &= 5(3x-2)^4(3) + 12(3x-2)^3(3) - 15 \\
 &= 15(3x-2)^4 + 36(3x-2)^3 - 15
 \end{aligned}$$

$$\begin{aligned}
 (3) \quad D_x[(x^2+3)^5 - 3(\frac{1}{x})^2 + 2(\sqrt{x})^3] \\
 &= 5(x^2+3)^4(2x) - 6(\frac{1}{x})(-\frac{1}{x^2}) + 6(\sqrt{x})^2 \cdot \frac{1}{2\sqrt{x}} \\
 &= 10x(x^2+3)^4 + \frac{6}{x^3} + 3\sqrt{x}
 \end{aligned}$$

Exercises 4-2c

1. (a) Prove Theorem 4-2c.
- (b) Prove Corollary 1 to Theorem 4-2c.
- (c) Prove Corollaries 2 and 3 to Theorem 4-2c.

2. Evaluate:

- | | |
|-----------------------------|-----------------------------------|
| (a) $D(x^5 - x^8 + x^{11})$ | (f) $D((3-5x)^2(1-x^2)^3)$ |
| (b) $D(5+x)^6$ | (g) $D(1-\frac{1}{x})^3$ |
| (c) $D(3-2x^3)^5$ | (h) $D(3x^{1/2} - 4x^{3/2})^6$ |
| (d) $D(x-3x^2+5x^3)^4$ | (i) $D(1-\sqrt{x})^{10}$ |
| (e) $D(x(1-x^2)^3)$ | (j) $D((3-5x+x^2)^3(1+x^2)^{10})$ |

3. Consider the curves $y = ax^3 + 1$ and $y = bx^2$.

- (a) Find two numbers a and b such that the curves have the same slope at $x = 1$ and that the sum of the slopes at $x = 2$ is 36.
- (b) Find values of a and b such that the curves have the same slope at a point of intersection.
- (c) Sketch the curves in part (b) for some allowable values of a and b .

4. (a) Sketch the curves $y = x^3$, $y = \frac{3}{x}$ for $|x| \leq 2$.
- (b) For what values of x is the derivative g' zero where $g(x) = x^3 + \frac{3}{x}$ for $|x| \leq 2$?
- (c) Find the discontinuities of g and g' .
- (d) Using the results of (a), (b), (c), sketch the graph of g for $-2 \leq x \leq 2$.
5. Let $u = f(x)$, $y = g(x)$, and $w = h(x)$.
- (a) Prove that if the functions f , g , and h are differentiable at x , then
- $$D(uvw) = (uv)Dh(x) + (uw)Dg(x) + (vw)Df(x).$$
- (b) Can you suggest a way to generalize your result to obtain a formula for the derivative of a product of n functions? Test your conjecture with the case $n = 4$.
- (c) Use the above result to evaluate:
- (i) $D[(5x - 2)(3 - 2x)(x^2 + 1)]$.
 - (ii) $D[(2x^3 - 3x^2 + 1)(\sqrt{x} + 1)^2]$.
 - (iii) $D[(3x - 2)(1 - x^2)(1 + x)(1 + x^2)]$.

Derivatives of quotients.

We have found the derivative $D\left(\frac{1}{x}\right)$ and, by the product rule, can obtain $D\left(\frac{1}{x^n}\right)$, but we still do not have a general rule for differentiating such functions as

$$f(x) = \frac{x - 1}{x^2 + 3}, \quad g(x) = \frac{\sqrt{x}}{x - 2}$$

or, more generally, the quotient of any two functions whose derivatives are known. Since the derivative of the general quotient $\frac{f(x)}{g(x)}$ can be obtained from that of the product

$$f(x) \cdot \left(\frac{1}{g(x)}\right),$$

we need only obtain the rule for the derivative of the reciprocal of $g(x)$.

THEOREM 4-2d. If $F(x) = \frac{1}{g(x)}$, then $F'(x) = \frac{-g'(x)}{(g(x))^2}$ at each point x for which $g'(x)$ exists and $g(x) \neq 0$.

$$\begin{aligned} \text{Proof. } \frac{F(x+h) - F(x)}{h} &= \frac{1}{h} \left[\frac{1}{g(x+h)} - \frac{1}{g(x)} \right] \\ &= \frac{1}{h} \left[\frac{g(x) - g(x+h)}{g(x)g(x+h)} \right] \\ &= \frac{-1}{g(x)g(x+h)} \left[\frac{g(x+h) - g(x)}{h} \right] \end{aligned}$$

From the theorems on limits and the continuity of $g(x)$ we obtain

$$F'(x) = \lim_{h \rightarrow 0} \frac{F(x+h) - F(x)}{h} = \frac{-g'(x)}{g(x)^2}.$$

Example 4-2f.

(1) If $f(x) = \frac{1}{x+2}$, then $f'(x) = \frac{-1}{(x+2)^2}$.

(2) If $g(x) = \frac{1}{x^n}$, then $g'(x) = \frac{-nx^{n-1}}{x^{2n}} = \frac{-n}{x^{n+1}}$.

Corollary 1. If f and g have derivatives at x and $g(x) \neq 0$, then for the quotient $G(x) = \frac{f(x)}{g(x)}$ we have

$$G'(x) = \frac{g(x)f'(x) - f(x)g'(x)}{(g(x))^2}.$$

The proof is left as an exercise.

Example 4-2g.

$$\begin{aligned} (1) \text{ If } G(x) = \frac{x-1}{x^2+3}, \text{ then } G'(x) &= \frac{(x^2+3)(1) - (x-1)(2x)}{(x^2+3)^2} \\ &= \frac{3+2x-x^2}{(x^2+3)^2}. \end{aligned}$$

- (2) We can sometimes simplify the differentiation of a rational function by giving it a convenient algebraic expression.

If $f(x) = \frac{x^3 + 3x^2 - 2}{x^2}$, then $f(x) = x + 3 - \frac{2}{x^2}$, and

$$f'(x) = 1 + \frac{4}{x^3}.$$

This technique is particularly helpful when the division can be performed rapidly.

- (3) Sometimes addition helps to shorten the work

$$g(x) = \frac{x+1}{x-1} - \frac{x}{x+1} = \frac{3x+1}{x^2-1}, \text{ and}$$

$$g'(x) = \frac{(x^2-1)(3) - (3x+1)(2x)}{(x^2-1)^2} = \frac{-3x^2 - 2x - 3}{(x^2-1)^2}.$$

Corollary 2. If R is a rational function, if $R(f(x))$ is defined, and if $f'(x)$ exists, then

$$D_x R(f(x)) = R'(f(x)) \cdot f'(x).$$

The proof is left as an exercise.

Exercises 4-2d

1. Evaluate:

(a) $D\left(\frac{x}{x-1}\right)$

(f) $D\left(\frac{\sqrt{x}}{1+x^2}\right)$

(b) $D\left(\frac{x^2}{1+x^2}\right)$

(g) $D\left(\frac{1}{1+\sqrt{x}}\right)$

(c) $D\left(1 - \frac{1}{x}\right)^{-1}$

(h) $D\left(\frac{x^2-1}{x^2+1}\right)^{-1}$

(d) $D\left(\frac{3+2x^2}{2-x^2}\right)$

(i) $D\left(\frac{1}{[\sqrt{x}]}\right)$

(e) $D\left(\frac{1}{x} + \frac{1}{1-x}\right)$

(j) $D\left(\frac{1}{\sqrt{[x]}}\right)$

2. Prove Corollary 1 to Theorem 4-2d.

3. Evaluate:

(a) $D\left(\frac{x^2 + 5}{x^2 - 5}\right)$

(f) $D\left(\frac{1 - x}{x^2(1 + x^2)}\right)$

(b) $D\left(\frac{x+1}{x} - \frac{x^2+1}{x^2}\right)$

(g) $D\left(\frac{1+x^{-1}}{1-x^{-1}}\right)$

(c) $D\left(\frac{x-1}{1+x+x^2}\right)$

(h) $D\left(\frac{ax^2 + bx + c}{dx^2 + ex + f}\right)$

(d) $D\left((1 + \frac{1}{x})(x+1)\right)$

(i) $D\left(\frac{x^4 + x^2 + 1}{x^2 + x + 1}\right)$

(e) $D\left(\frac{1-x+x^2}{1+x+x^2}\right)$

(j) $D\left(\frac{x^5 + x + 1}{x^2 + x + 1}\right)$

Note: These two reduce to particularly simple expressions.

4. Determine both

$$D\left(\frac{ax + b}{cx + d}\right) \text{ and } D\left(\frac{bc - ad}{c(cx + d)}\right).$$

Explain why both are the same.

5. Find the derivative of each of the following functions in as many ways as you can and describe its domain. (Do not overlook the definition of derivative.)

(a) $f : x \rightarrow \frac{|x|}{x}$

(c) $f : x \rightarrow \frac{x - [x]}{|x|}$

(b) $f : x \rightarrow \frac{x}{[x]}$

(d) $f : x \rightarrow \frac{(x - [x])^2}{|x|}$

6. Prove Corollary 2 to Theorem 4-2d.

7. Evaluate:

(a) $D_x \left(\frac{(x^2 - x^3)^2 - (x^2 - x^3)^4}{(x^2 - x^3)^2 + 1} \right)^7$

(b) $D_x \left(\frac{5(1 - \sqrt{x}) + 3}{(1 - \sqrt{x})^2 - 3} \right)^2$

(c) $D_x \left(\frac{(\sqrt{x} - \frac{1}{x})^4 - 1}{(\sqrt{x} - \frac{1}{x})^2 + 1} \right)^3$

8. Consider the quotient $\phi(x) = \frac{f(x)}{g(x)}$ where f and g have derivatives at x and where $g(x) \neq 0$. Obtain the formula for the derivative of a quotient by applying the product rule (Theorem 4-2b) to the expression

$$\phi(x)g(x) = f(x).$$

Why does this not constitute a proof of the quotient rule (Corollary 1 to Theorem 4-2d)?

4-3. Inverse Functions, Fractional Powers.

The preceding development can not be applied directly to find the derivatives of such functions as $f: x \rightarrow \sqrt[n]{x}$. We recall (see Appendix A2-4) that the number $y = \sqrt[n]{x}$ is defined as the principal solution of the equation

$$y^n = x$$

(namely, as the only solution of the equation when n is odd, and as the only nonnegative solution when n is even and $x \geq 0$). A natural approach to a discussion of the function f is through the familiar and well-understood function

$$g: y \rightarrow y^n.$$

The function g is inverse* to f ; that is, it undoes the effect of f . Thus, if f maps a onto b , i.e.,

$$f: a \rightarrow b,$$

then g maps b onto a , i.e.,

$$g: b \rightarrow a.$$

In this section we shall show in general how to differentiate the inverse of a function whose derivative we know.

From the property of inverses cited above it is easy to appreciate the graphical relation between inverse functions. If the point (a,b) is on the graph of f , then the point (b,a) is on the graph of g , and conversely; i.e., $b = f(a)$ if, and only if, $a = g(b)$. Since the points (a,b) and (b,a) are located symmetrically with respect to the line $y = x$, we observe that the graph $y = g(x)$ is the mirror image, in the line $y = x$, of the graph $y = f(x)$, as shown in Figure 4-3a. If the direction of the graph $y = f(x)$ at (a,b) is given with respect to the horizontal by the angle θ , then θ also gives the direction of the graph $y = g(x)$ at (b,a) with respect to the vertical. It follows that the slope of $y = g(x)$ at (b,a) is

$$\tan\left(\frac{\pi}{2} - \theta\right) = \cot \theta = \frac{1}{\tan \theta}.$$

(Of course, $\frac{1}{\tan \theta}$ is not defined if $\theta = 0$.) Intuitively, then, if f and g are differentiable functions and $f'(x) \neq 0$, we must have

*The symbol f^{-1} is often used for the inverse of f in other texts.

(1)

$$g'(b) = \frac{1}{f'(a)} \cdot$$

From the figure it is intuitive that if f has a derivative at a , then g has a derivative at b given by (1).

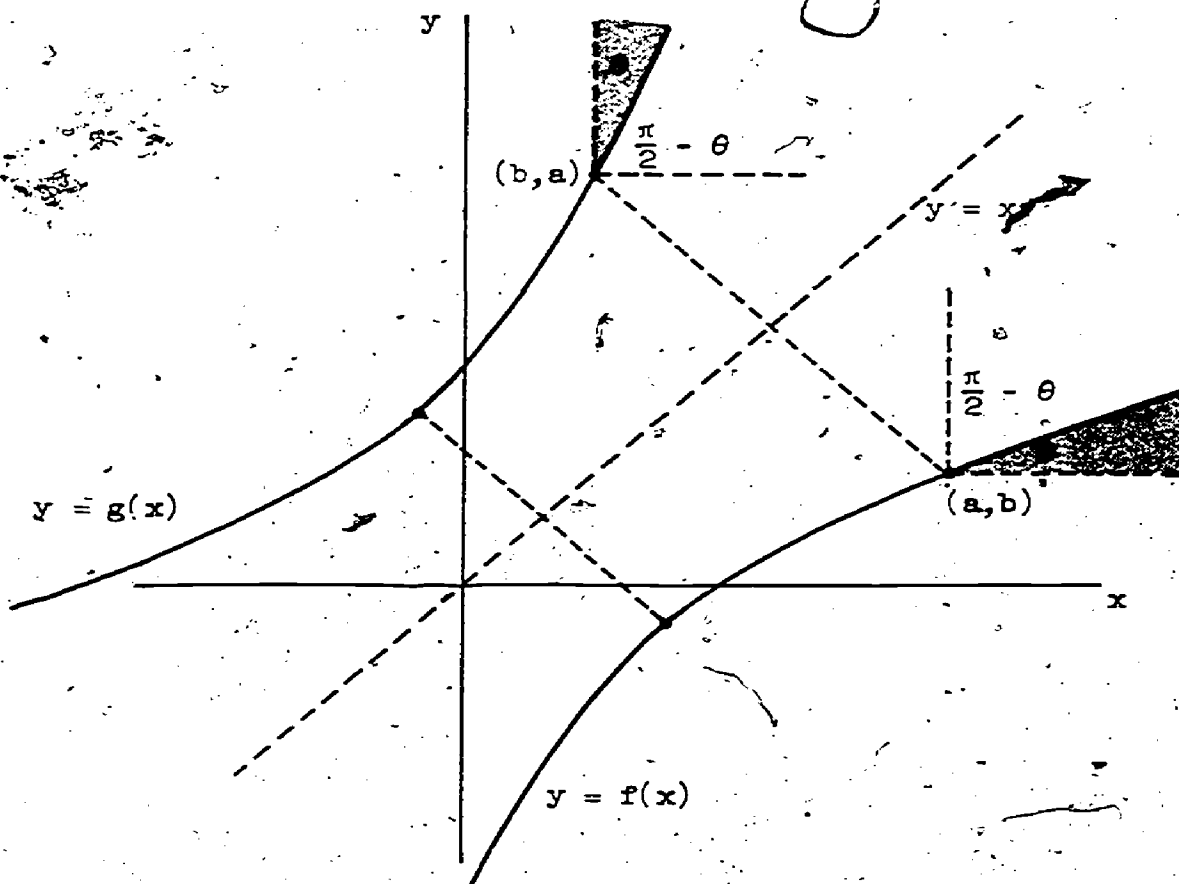


Figure 4-3a

We proceed to prove this formally as

THEOREM 4-3. Let f be either increasing or decreasing on a neighborhood of a . Then on that neighborhood f has an inverse g . If f has a derivative at a and $f'(a) \neq 0$, then g has a derivative at $b = f(a)$ and

$$g'(b) = \frac{1}{f'(a)} \cdot$$

Proof. On the neighborhood f is one-to-one and hence has a one-to-one inverse g , which is continuous at b by Theorem 3-6f. To investigate the existence and value of the derivative of g at b , we must consider the limit as y approaches b of

$$\begin{aligned}\frac{g(y) - g(b)}{y - b} &= \frac{g(y) - a}{f(g(y)) - f(a)} \\ &= \frac{1}{r(g(y))},\end{aligned}$$

where we have introduced the function

$$r(x) = \begin{cases} \frac{f(x) - f(a)}{x - a} & \text{for } x \neq a \\ f'(a) & \text{for } x = a \end{cases}$$

which is obviously continuous in the neighborhood of a . We apply in succession Theorem 3-6c and 3-6e on the continuity of quotients and compositions and obtain

$$\begin{aligned}g'(b) &= \lim_{y \rightarrow b} \frac{g(y) - g(b)}{y - b} = \lim_{y \rightarrow b} \frac{1}{r(g(y))} \\ &= \frac{1}{r(g(b))} = \frac{1}{r(a)} = \frac{1}{f'(a)}\end{aligned}$$

as desired.

Note that the theorem states that f' and g' are reciprocals at different points (see Figure 4-3a); f' at $x = a$ in the domain of f and g' at $y = b$ in the range of f . For example, let $f : x \rightarrow x^2 + 1$ for $x > 0$. Then $g : x \rightarrow \sqrt{x - 1}$ for $x > 1$. $f'(x) = 2x$, so

$$f'(g(x)) = 2\sqrt{x - 1}, \text{ and } g'(x) = \frac{1}{f'(g(x))} = \frac{1}{2\sqrt{x - 1}}.$$

Thus, for example

$$g'(10) = \frac{1}{f'(g(10))} = \frac{1}{f'(3)} = \frac{1}{6}; \quad g'(3) = \frac{1}{f'(g(3))} = \frac{1}{f'(\sqrt{2})} = \frac{1}{2\sqrt{2}}.$$

Exercises 4-3a

1. Show that $f : x \rightarrow \frac{1-x}{1+x}$, where $x > -1$, has an inverse. Find an equation that defines the inverse g and find the derivative of g .
2. Verify that the inverse of $f : x \rightarrow x + |x|$ exists and then find the derivative of the inverse.
3. Sketch the graph of $f : x \rightarrow x^3 - 3x$ and tell why f does not have an inverse. Indicate how you can divide the domain of f into three parts, and define three new functions each of which agrees with f on its domain and has an inverse. Justify your result. (See Exercises 2-3, No. 6.)
4. Consider the function defined by

$$f(x) = \begin{cases} 2x & , \text{ for } x > 0 \text{ and irrational} \\ x^2 + 1 & , \text{ for } x \geq 0 \text{ and rational} \end{cases}$$

Show that f has a derivative at $x = 1$. Prove that the mapping f is one-to-one, i.e., that f has an inverse. Prove that the inverse of f is not differentiable at any point.

As a first and important application of Theorem 4-3, we compute the derivative of the n -th root function. The function

$$g : y \rightarrow \sqrt[n]{y} = y^{1/n}$$

is defined as the inverse of the n -th power function

$$f : x \rightarrow x^n.$$

Here n is any natural number. We restrict the domain of f to nonnegative numbers (Why? See Exercise 4-3b, No. 4). With $b = a^n$, $a > 0$, we find

$$g'(b) = \frac{1}{f'(a)} = \frac{1}{n a^{n-1}} = \frac{1}{n} a^{1-n} = \frac{1}{n} (b^{1/n})^{1-n} = \frac{1}{n} b^{\frac{1}{n} - 1}$$

That is, the formula $Dx^r = rx^{r-1}$, previously established for integers r , holds also for $r = \frac{1}{n}$.

We can establish the same formula for any rational number $r = \frac{p}{q}$, p and q integers, $q \geq 1$. Using Theorem 4-2c, we have

$$Dx^r = D(x^{1/q})^p = p(x^{1/q})^{p-1} D(x^{1/q}) = p(x^{1/q})^{p-1} \cdot \frac{1}{q} x^{\frac{1}{q}-1} = \frac{p}{q} x^{\frac{p}{q}-1} = rx^{r-1}.$$

Thus we have

Corollary. For every rational number r , $Dx^r = rx^{r-1}$ ($x > 0$).

The extension of this corollary to general real (including irrational) powers is deferred until Chapter 8.

Exercises 4-3b

1. Evaluate the following and express your answers using positive exponents only.

(a) $Dx^{2/3}$

(b) $Dx^{3/5}$

(c) $Dx^{-2/3}$

(d) $Dx^{-3/5}$

2. Find $f'(2)$ if:

(a) $f(x) = (2x)^{1/3}$

(b) $f(x) = x^{-1/4}$

(c) $f(x) = x^{-4/5}$

3. Evaluate:

(a) $D(x^{1/3} - 3x^{-1/3})$

(b) $D\left(\frac{1 + x^{4/3}}{1 - x^{4/3}}\right)$

(c) $D\left(\sqrt{x^3} - \sqrt{x} + \frac{1}{\sqrt{x}} - \frac{1}{\sqrt{x^3}}\right)$

4. Let $f(x) = x^n$ for n an integer. For what points a in the domain of f do the hypotheses of Theorem 4-3 hold?
5. Consider the function $f : x \rightarrow \sqrt[n]{x}$. We have $f(x)^n = x$. Applying Theorem 4-2c, obtain $D(\sqrt[n]{x}) = \frac{1}{n} x^{\frac{1}{n} - 1}$. Why is this not a proof of the corollary to Theorem 4-3?
6. Let $r = \frac{p}{q}$ under the hypothesis of the corollary to Theorem 4-3 and let q be odd. Prove the corollary for $x < 0$. Why is the case $x = 0$ not included?
7. Under the conditions of the preceding exercise, show that for $g : x \rightarrow x^{p/q}$, where $p > q$, the derivative of g at zero exists and $g'(0) = 0$.

4-4. Circular Functions.

In Section 2-5 we ambitiously attacked the problem of evaluating $D \sin x$ at $x = 0$. We reduced the problem to that of evaluating

$$(1) \quad \alpha = \lim_{h \rightarrow 0} \frac{\sin h}{h}$$

under the assumption that the limit exists. As we shall now see, the evaluation of $D \sin x$ at any point can be reduced to the evaluation of this limit. Using the formula for the sine of a sum, we have

$$\frac{\sin(x+h) - \sin x}{h} = \frac{\sin x \cos h + \cos x \sin h - \sin x}{h} = \frac{\sin h}{h} \cos x + \frac{\cos h - 1}{h} \sin x.$$

From the theorems on limits,

$$D \sin x = \alpha \cos x + \beta \sin x$$

where $\beta = \lim_{h \rightarrow 0} \frac{\cos h - 1}{h}$. We shall assume for the present that α defined in (1) exists, and from this will prove that $\beta = 0$.

Since $\cos h - 1 = -2 \sin^2 \frac{h}{2}$, we obtain

$$\frac{\cos h - 1}{h} = - \left(\frac{\sin \frac{h}{2}}{\frac{h}{2}} \right)^2 \left(\frac{h}{2} \right)^{-1}.$$

Since $\lim_{h \rightarrow 0} \frac{\sin \frac{h}{2}}{\frac{h}{2}} = \lim_{k \rightarrow 0} \frac{\sin k}{k} = \alpha$, and $\lim_{h \rightarrow 0} \frac{h}{2} = 0$ we get by using theorems

on limits of products of functions that $\beta = 0$. Consequently

$$(2) \quad D \sin x = \alpha \cos x$$

where the constant of proportionality α has yet to be determined.

In reviewing the preceding argument we find several matters assumed without proof:

(a) The formula for the sine of a sum:

$$\sin(u + v) = \sin u \cos v + \cos u \sin v.$$

(b) $\sin 0 = 0$

(c) $1 - \cos u = 2 \sin^2 \frac{u}{2}$.

(d) the existence of $\lim_{h \rightarrow 0} \frac{\sin h}{h}$.

It may seem odd that the well-known properties (a) - (c) of the sine and cosine are listed as assumptions. The rules of Section A2-5 derive from an intuitively based idea of length for circular arcs. The concept of length for curves other than polygons is defined in a later chapter as a limit of approximations by polygonal arcs. Clearly, our knowledge of the circular functions does not rest upon any such precise definition but upon geometrical intuition. Similarly, we derive property (d) not by analytical reasoning but by arguing plausibly from a picture. Later we shall see how the argument can be made analytically complete (Section 8-5), but the circular functions are too important for us to delay our account until the gaps can be filled. In fact, we shall profit by gaining an intuitive understanding of the circular functions before attempting to be formally precise.

Consider a ray from the origin lying in the first quadrant. (Figure 4-4).

If, on the unit circle, x is the length of the arc between the ray and the positive horizontal axis, then the ray intersects the circle at the point with coordinates $(\cos x, \sin x)$.

There are two similar triangles in the figure: the smaller has base $\cos x$ and altitude $\sin x$; the larger has base 1 and altitude $\tan x$

(since $\frac{\sin x}{\cos x} = \frac{\tan x}{1}$). Since the

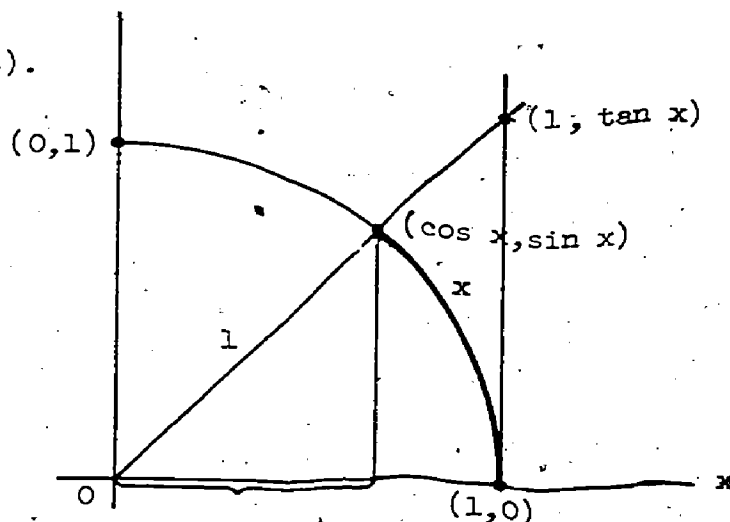


Figure 4-4

circular sector determined by the arc contains one triangle and is contained in the other, its area, $\frac{x}{2}$, lies between their areas:

$$\frac{1}{2} \cos x \sin x < \frac{x}{2} < \frac{1}{2} \cdot 1 \cdot \tan x.$$

Here we have used the intuitively evident fact that the area of a sector is proportional to the length of the corresponding arc.

The factor of proportionality is $\frac{1}{2}$, since the complete unit circle has area π and circumference 2π . Later it will be easy to obtain the constant of proportionality analytically.

On multiplying by the positive value, $\frac{2}{\sin x}$, in the preceding inequality, we obtain

$$\cos x < \frac{x}{\sin x} < \frac{1}{\cos x},$$

whence

$$(3) \quad \frac{1}{\cos x} > \frac{\sin x}{x} > \cos x.$$

Since $\sin(-x) = -\sin x$ and $\cos(-x) = \cos x$, this inequality is also valid for x in the fourth quadrant, $-\frac{\pi}{2} < x < 0$. Since $\cos \theta$ is continuous at $\theta = 0$, we have*

$$\lim_{x \rightarrow 0} \cos x = \cos 0 = 1,$$

and $\lim_{x \rightarrow 0} \frac{1}{\cos x} = 1$. This is precisely the kind of situation in which we can apply the Squeeze Theorem (Corollary 2 to Theorem 3-4f). It follows immediately from (3) that

$$\lim_{x \rightarrow 0} \frac{\sin x}{x} = 1.$$

Entering this result ($\alpha = 1$) in Equation (2) we have, finally,

$$D \sin x = \cos x.$$

The proof of the analogous formula,

$$D \cos x = -\sin x,$$

is left as an exercise.

The derivatives of the other circular functions are now easily obtained. For instance, from the theorem on the differentiation of a quotient, we have

$$\begin{aligned} D \tan x &= D \frac{\sin x}{\cos x} = \frac{\cos x D \sin x - \sin x D \cos x}{(\cos x)^2} \\ &= \frac{\cos^2 x + \sin^2 x}{\cos^2 x} = \frac{1}{\cos^2 x} = \sec^2 x = 1 + \tan^2 x. \end{aligned}$$

We leave the problem of differentiating the other circular functions as an exercise, but list some of the results for easy reference:

- (a) $D \sin x = \cos x$
- (b) $D \cos x = -\sin x$
- (c) $D \tan x = 1 + \tan^2 x$
- (d) $D \cot x = -1 - \cot^2 x$

* Condition (d) above implies the continuity of $\sin \theta$ at $\theta = 0$ and the continuity of $\cos \theta$ at $\theta = 0$ then follows from (c).

Exercises 4-4

1. Show that $D \cos x = -\sin x$.
2. Evaluate $\lim_{h \rightarrow 0} \frac{\tan h}{h}$. (Hint: Express $\tan h$ in terms of $\sin h$ and $\cos h$.)
3. From the definition of the derivative as a limit and the result of No. 2 derive the formula

$$D \tan x = \sec^2 x.$$

[Hint: $\tan(x + h) = \frac{\tan x + \tan h}{1 - \tan x \tan h}$.]

Compare this result with the result obtained by using the method of differentiating the quotient $\frac{\sin x}{\cos x}$.

4. In the simplest way you can, evaluate the following and express your answers in several different equivalent forms.

- (a) $D \cot x$
- (b) $D \sec x$
- (c) $D \csc x$
- (d) $D \sin^2 x$
- (e) $D \cos^2 x$
- (f) $D(4 \cos^3 x - 3 \cos x)$
- (g) $D(3 \sin x - 4 \sin^3 x)$

5. Evaluate the following limits.

- (a) $\lim_{h \rightarrow 0} \frac{\sin 2h}{h}$
- (b) $\lim_{h \rightarrow 0} \frac{1 - \cos h}{h^2}$

6. (a) Given that $\lim_{x \rightarrow 0} \sin x = \sin 0 = 0$.

Prove that $\lim_{x \rightarrow 0} \frac{\sin x}{\cos x + 1} = 0$. (Hint: Show that $\cos x$ is continuous at $x = 0$.)

- (b) From the preceding result prove that $\sin x$ and $\cos x$ are continuous for all values of x .

Make explicit just what is being assumed in the proofs of (a) and (b).

7. Given that $Df(x) = G(x)$, show that $Df(ax + b) = aG(ax + b)$, provided f is differentiable at $ax + b$.
8. Evaluate the following.
- $D(\cos^2 x)(\sin 2x)$
 - $D \sin^2(ax + b)$
 - $D(\sin 7x)(\cos 2x)$
9. Let $g(x) = |\cos x|$. Discuss the domain of the derivative for x in the interval $0 \leq x \leq \pi$.
10. Find a point on the graph of $y = \sin x$ at which the slope of the curve is equal to the slope of the line $x + 2y + 2 = 0$. Is there only one such point? Justify your answer.
11. Evaluate the following.
- $D \left(\frac{1 - \sin x}{1 + \cos^2 x} \right)$
 - $D \left(\frac{1 - \tan^2 x}{2 \tan x} \right)$
 - $D \left(\frac{\sin^4 x}{\cos^2 x} \right)$
 - $D \left(\frac{3}{\cos^2 x} \right)$
 - $D \left(\frac{1}{1 + \tan x} \right)$
 - $D(x \tan x)$
 - $D \left(\frac{\sin x + \cos x}{\sin x - \cos x} \right)^2$
12. Show that there are no points of the graph of $y = \sec x - \tan x$ at which the slope of the curve is zero.
13. Find all values of x for which the slope of the graph of $f(x) = \sin x \tan x$ is zero.

14. (a) Sketch the graph of $f(x) = \frac{x}{\sin x}$ for $0 < x < \frac{\pi}{2}$.
- (b) Examine $f'(x)$ and show that there is no value of x in the interval $0 < x < \frac{\pi}{2}$ for which $f'(x) = 0$.
- (c) Explain how your results support the fact that f is increasing on the given interval.

15. (a) Find the maximum and minimum values of the function

$$f(x) = a \cos x + b \sin x.$$

- (b) Sketch the graph of f .

16. Consider $f : x \rightarrow \sin \frac{1}{x}$ in the domain $0 < x \leq 1$. Is it possible to define f at $x = 0$ such that the function is continuous in $[0, 1]$?

17. Consider $f : x \rightarrow x \sin \frac{1}{x}$ for $x \neq 0$.

- (a) Sketch the graph of f .

- (b) Is it possible to define f at $x = 0$ such that the function is continuous at $x = 0$?

- A18. Does the function

$$f : x \rightarrow \begin{cases} x^2 \sin \frac{\pi}{x}, & x \neq 0 \\ 0, & x = 0, \end{cases}$$

have a derivative at $x = 0$?

19. Given that the functions S and C satisfy the equations

$$DS = C,$$

$$DC = -S.$$

Show that $D(S^2 + C^2) = 0$.

20. Given that the functions f , g , and h satisfy the equations

$$Df = g,$$

$$Dg = h,$$

$$Dh = f.$$

Show that $D(f^3 + g^3 + h^3 - 3fgh) = 0$.

4-5. Inverse Circular Functions.

The most striking feature of the graphs of the circular functions is their periodic or cyclic aspect, the uniform iteration of the same geometric pattern. For the sine and cosine, this periodic property is expressed analytically by the relations

$$\begin{aligned}\sin(x + 2\pi) &= \sin x \\ \cos(x + 2\pi) &= \cos x,\end{aligned}$$

where 2π is seen to be a period of these functions. From the representation of the other circular functions in terms of the sine and cosine it follows that they, too, have the period 2π . However, the tangent and cotangent have the shorter period π as well.

From the periodic character of any circular function f , it follows at once that f cannot represent a one-to-one mapping on its entire domain: if $b = f(a)$, then $b = f(a + 2k\pi)$ for all integers k . Every point of the range of f is covered infinitely many times in the mapping. We cannot obtain an inverse for f over its entire domain, so we restrict ourselves to a special representative interval where either f is increasing throughout the interval or f is decreasing throughout the interval. The names of the inverses of the sine, cosine, etc., on such restricted domains, are arcsine, arccosine, etc.*

For the sine function we choose the representative interval $[-\frac{\pi}{2}, \frac{\pi}{2}]$ where the sine is increasing (see Figure 4-5a). The range of the sine on this interval is the interval $[-1, 1]$. The arcsine, therefore, has the domain $[-1, 1]$ and the range $[-\frac{\pi}{2}, \frac{\pi}{2}]$. Consequently we define the value $y = \arcsin x$ in the mapping

$$\arcsin : x \longrightarrow y, \quad \text{for } -1 \leq x \leq 1$$

as the one value satisfying $-\frac{\pi}{2} \leq y \leq \frac{\pi}{2}$ for which $\sin y = x$.

We could have used any other representative interval, $[k\pi - \frac{\pi}{2}, k\pi + \frac{\pi}{2}]$, for $k = 0, \pm 1, \pm 2, \dots$, to define an inverse of the sine. For this reason the specific inverse, the arcsine with range $[-\frac{\pi}{2}, \frac{\pi}{2}]$, is known as

* A common alternative notation is \sin^{-1} for arcsin, \tan^{-1} for arctan, and so on. We do not use it in this text because it conflicts with the universally accepted convention of writing $\sin^n x$ for $(\sin x)^n$, etc.

the principal inverse of the sine. Similarly the specific inverse circular functions defined below are known as principal inverses of their respective circular functions.

For the cosine we choose the representative interval $[0, \pi]$ where the cosine is decreasing (see Figure 4-5b). We define the value $y = \arccos x$ for $-1 \leq x \leq 1$ as the one value satisfying $0 \leq y \leq \pi$ for which $\cos y = x$.

Finally, for the tangent we use the representative interval $(-\frac{\pi}{2}, \frac{\pi}{2})$ and define $y = \arctan x$ for any x as the value in the interval $-\frac{\pi}{2} < y < \frac{\pi}{2}$ for which $\tan y = x$ (see Figure 4-5c).

We need not concern ourselves (in the text) with inverse functions for the cotangent, secant, and cosecant; these can be treated in terms of the functions already at our disposal (see Exercises 4-5) and are used infrequently.

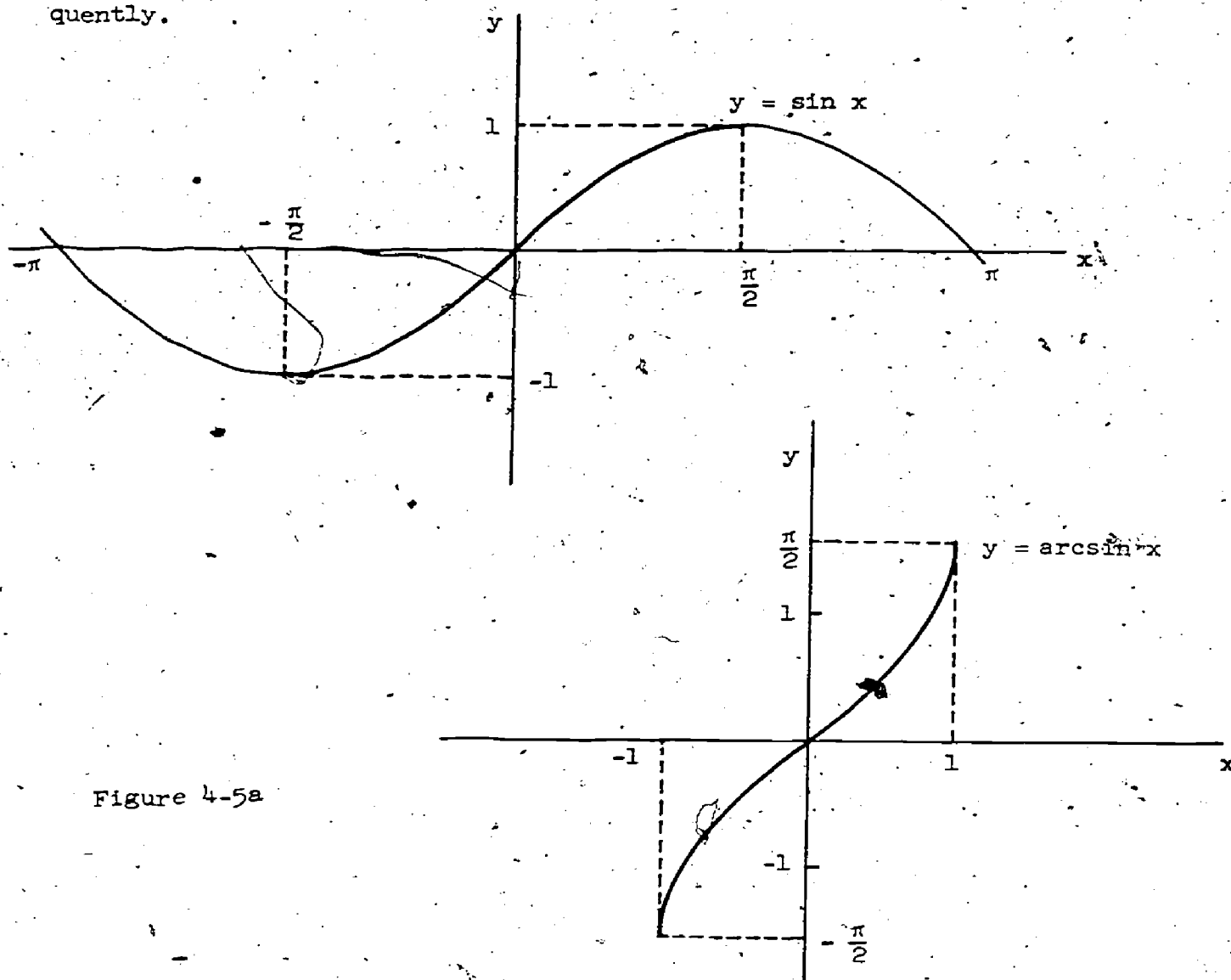


Figure 4-5a

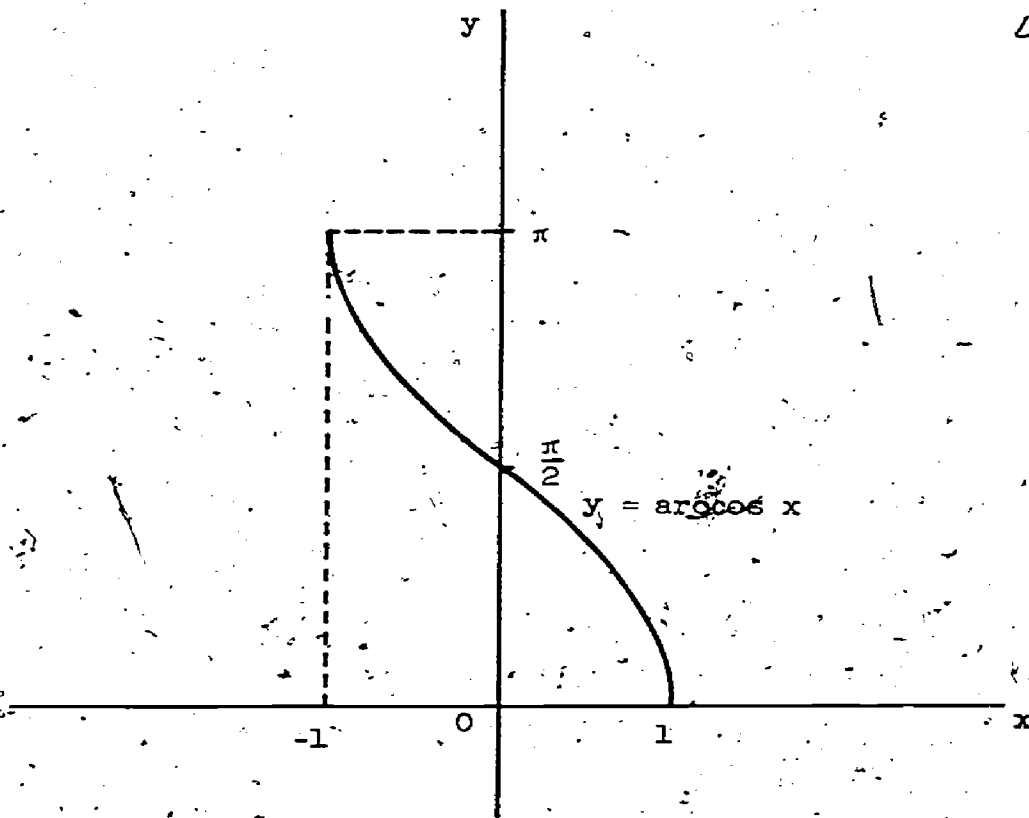
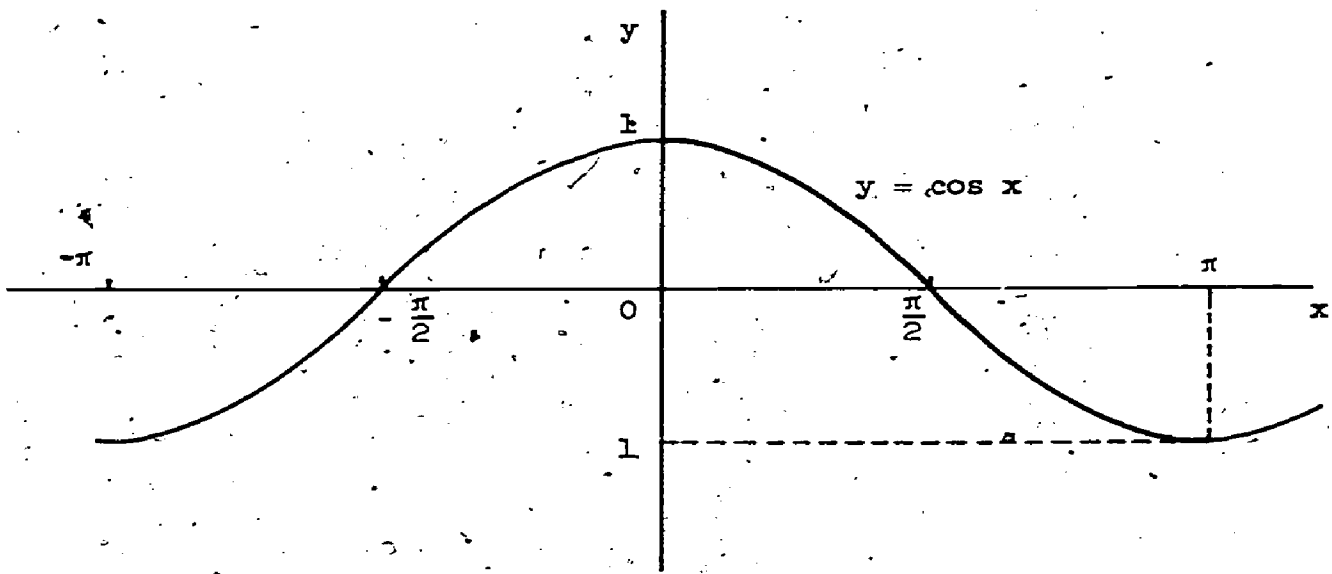


Figure 4-5b

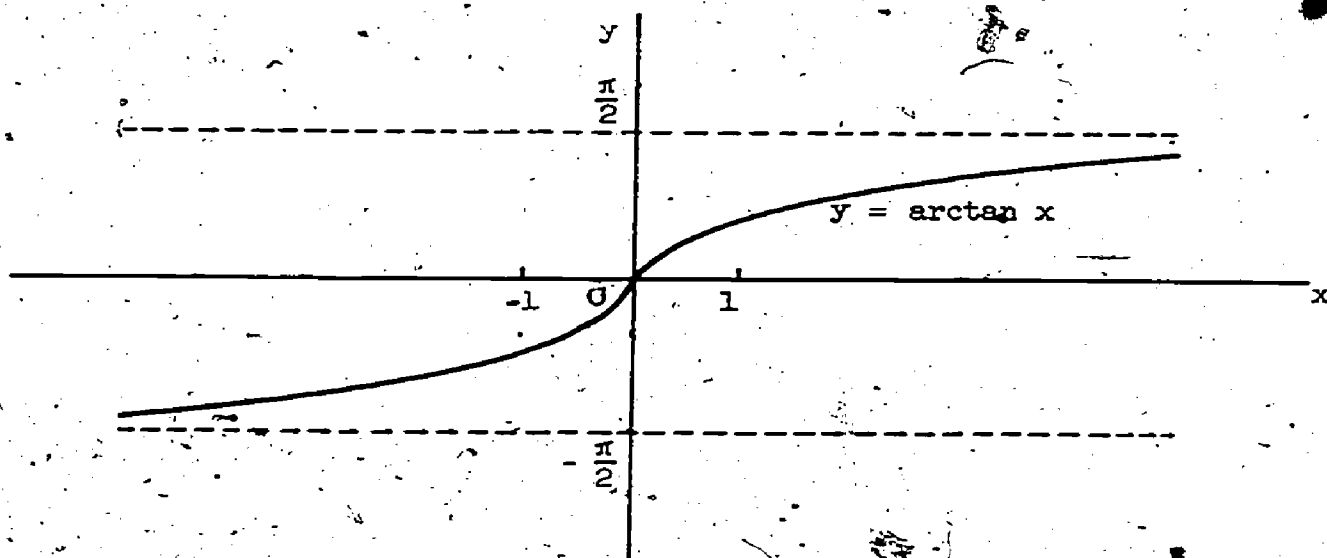
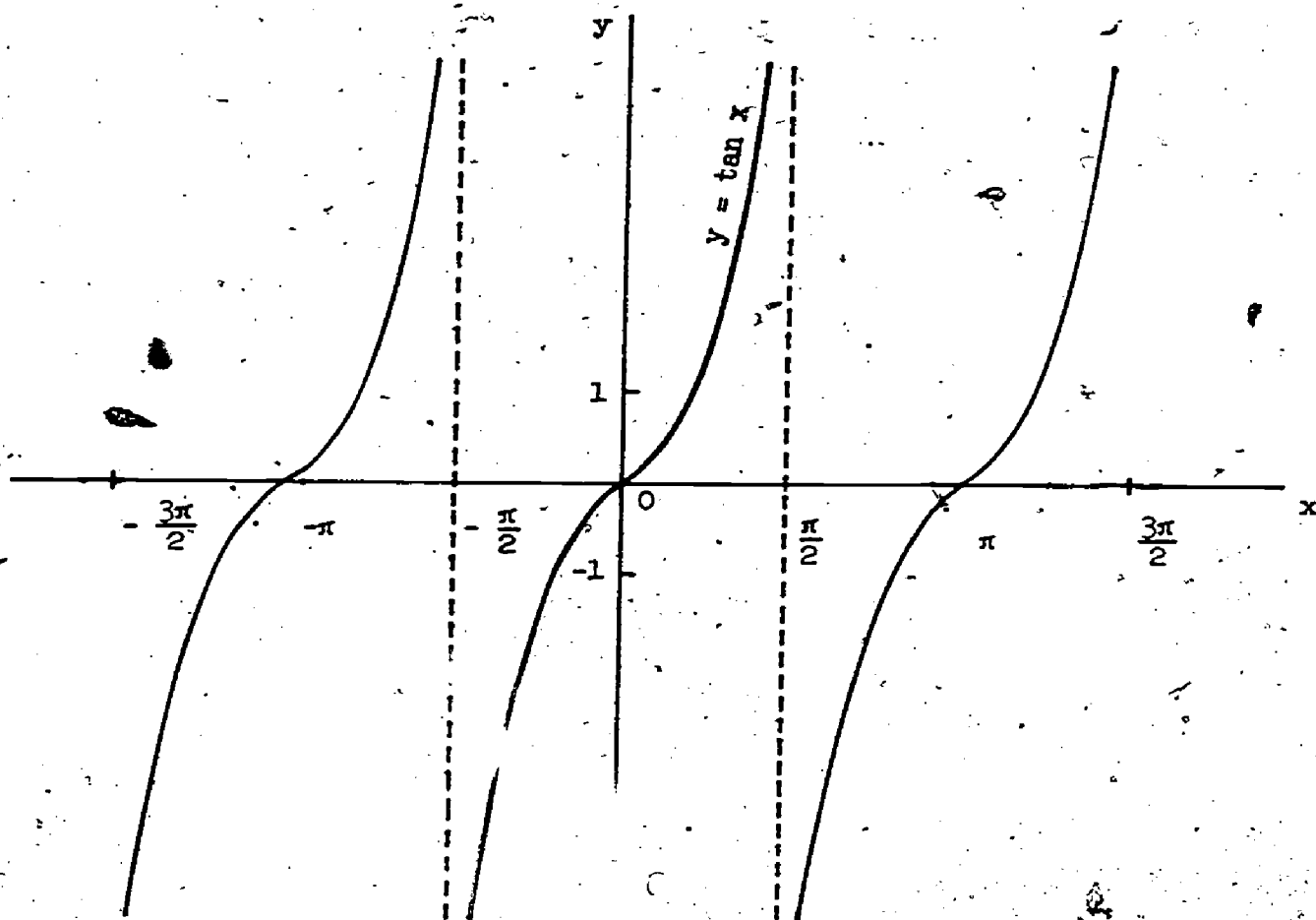


Figure 4-5c

We now turn to the task of finding the derivatives of the inverse circular functions. This is merely a matter of applying the general theorem on derivatives of inverses. The function

$$g : y \longrightarrow \arcsin y, \quad \text{for } -1 \leq y \leq 1,$$

is the inverse of

$$f : x \longrightarrow \sin x, \quad \text{for } -\frac{\pi}{2} \leq x \leq \frac{\pi}{2},$$

with derivative $f'(x) = \cos x$. Hence if $b = \sin a$, $-\frac{\pi}{2} < a < \frac{\pi}{2}$, then

$$g'(b) = \frac{1}{f'(a)} = \frac{1}{\cos a} = \frac{1}{\sqrt{1 - \sin^2 a}} = \frac{1}{\sqrt{1 - b^2}}.$$

Therefore

$$D \arcsin x = \frac{1}{\sqrt{1 - x^2}}, \quad \text{for } -1 < x < 1.$$

As an exercise you are asked to derive in the same way the formulas

$$D \arccos x = \frac{-1}{\sqrt{1 - x^2}}, \quad \text{for } -1 < x < 1,$$

$$D \arctan x = \frac{1}{1 + x^2}, \quad \text{for all } x.$$

(The first of these follows immediately from the identity $\arccos x = \frac{\pi}{2} - \arcsin x$.) It is noteworthy that the derivatives of all the inverse circular functions are algebraic functions:

Exercises 4-5

1. Determine the domain and range and draw the graph of the function

(a) $f : x \longrightarrow \arcsin(\sin x)$.

(b) $f : x \longrightarrow \sin(\arcsin x)$.

(c) $f : x \longrightarrow \arcsin(\cos x)$.

(d) $f : x \longrightarrow \cos(\arcsin x)$.

(e) $f : x \longrightarrow \arctan(\tan x)$.

2. Derive the formula

$$D \arccos x = \frac{-1}{\sqrt{1-x^2}},$$

without using the identity which relates $\arccos x$ to $\arcsin x$.

3. Derive the formula

$$D \arctan x = \frac{1}{1+x^2}.$$

4. Derive each of the following formulas.

$$(a) \quad D \operatorname{arccot} x = -\frac{1}{1+x^2}$$

$$(b) \quad D \operatorname{arcsec} x = \frac{1}{|x|\sqrt{x^2-1}}$$

$$(c) \quad D \operatorname{arccsc} x = \frac{-1}{|x|\sqrt{x^2-1}}$$

5. Evaluate:

$$(a) \quad D (\arcsin x + \arccos x)$$

$$(b) \quad D (x^2 \arcsin x)$$

$$(c) \quad D \frac{x^2}{\arctan x}$$

$$(d) \quad D (\arcsin x)^3$$

$$(e) \quad D \frac{1}{1 + \arcsin x}$$

6. Find $\lim_{h \rightarrow 0} \frac{\arcsin h}{h}$. (Hint: What is the definition of the derivative of $f(x) = \arcsin x$ at $x = 0$?)

7. Evaluate:

$$(a) \quad D \left(\frac{\arcsin x}{1 + \arccos x} \right)$$

$$(b) \quad D \left(\frac{\operatorname{arcsec} x}{1 - \operatorname{arccsc} x} \right)$$

$$(c) \quad D \left(\frac{1 - \arctan x}{1 + \arctan x} \right)$$

4-6. Compositions. Chain Rule.

We have seen that we may construct new functions by composition of already known functions. As we shall now show, the derivative of a composition can be expressed in terms of its constituent functions.

Consider the function $x \longrightarrow gf(x)$ obtained by composition of the functions g and f . If g is a polynomial, then we have the formula

$$(1) \quad D_x gf(x) = g'(f(x)) D_x f(x) = g'(f(x)) f'(x)$$

(Corollary 3 to Theorem 4-2c); but we have not derived a corresponding formula if g is not a polynomial.

The formula (1) is an instance of a general principle called the Chain Rule because it can be used to differentiate a chain of compositions, for example

$$(2) \quad D_x hgf(x) = h'(gf(x)) D_x gf(x) = h'(gf(x)) g'(f(x)) f'(x).$$

THEOREM 4-6. Let the function f be differentiable at a , and let the function g be differentiable at $b = f(a)$. Then the composition $\phi : x \longrightarrow gf(x)$ has a derivative at a given by

$$\phi'(a) = g'(f(a)) f'(a).$$

Before proceeding with the actual proof let us see first how the result stated by the theorem might have been conjectured. We may view the definition for the derivative at a of a function: $x \longrightarrow h(x)$, $\lim_{x \rightarrow a} \frac{h(x) - h(a)}{x - a}$, as follows. Call $x - a = \Delta a$ the "increment of x at a " and $h(x) - h(a) = \Delta h(a)$ the "induced increment of $h(x)$ at a ." Then $h'(a)$ may be interpreted as the limit of ratios of increments $\frac{\Delta h(a)}{\Delta a}$ as Δa approaches zero. Using this interpretation we would anticipate that

$$\phi'(a) = \lim_{\Delta a \rightarrow 0} \frac{\Delta \phi(a)}{\Delta a} = \lim_{\Delta a \rightarrow 0} \frac{\Delta gf(a)}{\Delta f(a)} \cdot \frac{\Delta f(a)}{\Delta a}.$$

This together with the theorem on limits of products and compositions suggests that the anticipated formula for $\phi'(a)$ is the correct one. Our deduction would lead to a proof of the theorem if $\Delta f(a)$ were never zero. This cannot be guaranteed in all cases; for example, if f is a constant function, then $\Delta f(a)$ is always zero. Thus we must proceed with greater care, taking note of situations in which $\Delta f(a)$ may be zero at some points.

Proof. If $f(x) \neq f(a) = b$, then

$$\frac{\phi(x) - \phi(a)}{x - a} = \frac{gf(x) - gf(a)}{x - a} = \frac{gf(x) - g(b)}{f(x) - b} \cdot \frac{f(x) - f(a)}{x - a}.$$

If $f(x) = b$, it is true that for any number c

$$\frac{\phi(x) - \phi(a)}{x - a} = c \frac{f(x) - f(a)}{x - a}$$

since both sides are zero. In terms of the function r defined by

$$r(y) = \begin{cases} \frac{g(y) - g(b)}{y - b}, & \text{for } y \neq b, \\ c, & \text{for } y = b, \end{cases}$$

we then have

$$(3) \quad \frac{\phi(x) - \phi(a)}{x - a} = r(f(x)) \frac{f(x) - f(a)}{x - a}.$$

We could apply the theorem on limits if r were continuous at b . It is, however, not continuous unless we choose c judiciously. Pick $c = g'(b)$.

Then $\lim_{y \rightarrow b} r(y) = \lim_{y \rightarrow b} \frac{g(y) - g(b)}{y - b} = g'(b) = r(b)$, that is, r is continuous

at b . From (3) and Theorem 3-6e on the continuity of compositions, we obtain now

$$\begin{aligned} \phi'(a) &= \lim_{x \rightarrow a} \frac{\phi(x) - \phi(a)}{x - a} = \lim_{x \rightarrow a} r(f(x)) \lim_{x \rightarrow a} \frac{f(x) - f(a)}{x - a} \\ &= r(f(a)) f'(a) = g'(f(a)) f'(a). \end{aligned}$$

Example 4-6a. Evaluate $D \sin x^2$. With $f(x) = x^2$ and $g(y) = \sin y$, we have $f'(x) = 2x$, $g'(y) = \cos y$, and therefore,

$$D \sin x^2 = 2x \cos x^2.$$

Example 4-6b. Differentiate $\phi(x) = \arctan \frac{1}{x}$. We set $f(x) = \frac{1}{x}$ and $g(y) = \arctan y$, so that

$$f'(x) = \frac{-1}{x^2}, \quad g'(y) = \frac{1}{1 + y^2}, \quad \text{and}$$

$$\phi'(x) = D \arctan \frac{1}{x} = \frac{1}{1 + (\frac{1}{x})^2} \left(\frac{-1}{x^2} \right) = \frac{-1}{x^2 + 1}$$

Example 4-6c. Evaluate $D \cos\left(\frac{1-x}{1+x}\right)^{2/3}$. We use (2), with

$$f(x) = \frac{1-x}{1+x}, \quad g(y) = y^{2/3}, \quad h(z) = \cos z,$$

$$f'(x) = \frac{-2}{(1+x)^2}, \quad g'(y) = \frac{2}{3} y^{-1/3}, \quad h'(z) = -\sin z.$$

The formula yields

$$\begin{aligned} D \cos\left(\frac{1-x}{1+x}\right)^{2/3} &= -\sin\left(\frac{1-x}{1+x}\right)^{2/3} \cdot \frac{2}{3} \left(\frac{1-x}{1+x}\right)^{-1/3} \cdot \frac{-2}{(1+x)^2} \\ &= \frac{4}{3} (1-x)^{-1/3} (1+x)^{-5/3} \sin\left(\frac{1-x}{1+x}\right)^{2/3}. \end{aligned}$$

Exercises 4-6

1. For each of the following find $Dfg(x)$, $Dgf(x)$, $Dff(x)$, and $Dgg(x)$.

(a) $f(x) = x^3 - 2x$, $g(x) = \sqrt{x}$

(b) $f(x) = \sin x$, $g(x) = \cos x$

(c) $f(x) = x^2$, $g(x) = \sin x$

(d) $f(x) = \frac{x}{1+x}$, $g(x) = \frac{x^2}{1+x}$

(e) $f(x) = \sin(x^2)$, $g(x) = \sqrt{1-x^2}$

2. Find $Dff(x)$ for:

(a) $f(x) = \sin^2 x$.

(b) $f(x) = \tan^2 x$.

3. Find $f'(x)$ if $f(x)$ is:

(a) $(x^3 + 4)^{1/2}$

(b) $(2x^2 + 2)^{-1/2}$

(c) $\frac{\sqrt{1-3x} + 1}{\sqrt{1-3x}}$

(d) $\sqrt{\sin^2 x + x^2}$

4. Evaluate:

(a) $D_x \sqrt{a^2 - x^2}$

(b) $D_x \left((x^2 + 1)^{1/2} + (x^2 + 1)^{-1/2} \right)$

(c) $D_x \left(\frac{\sqrt{x^2 - a^2}}{\sqrt{x^2 + a^2}} \right)$

(d) $D_x \left(\frac{1 + \sqrt{1 - 2x}}{\sqrt{1 - 2x}} - \sqrt{1 - 2x} \right)$

(e) $D_x \left(x(2x^2 - 2x + 1)^{-1/2} \right)$

5. Evaluate:

(a) $D_x (\sqrt{1 + \cos x})$

(e) $D_x (\arctan (\arctan x))$

(b) $D_x (x^2 \sqrt{\sin x})$

(f) $D_x (x^2 \sin x \cos x)$

(c) $D_x (\cos (\cos (\cos x)))$

(g) $D_x \left(\frac{\sin^2 x}{\sin(x^2)} \right)$

(d) $D_x (\arcsin(\cos x))$

(h) $D_x \left(\tan \left(\frac{1+x}{1-x} \right) \right)$

6. Evaluate:

(a) $D_x (\arcsin(\sin x - \cos x))$

(f) $D_x (\operatorname{arcsec} \sqrt{1 + x^2})$

(b) $D_x \left(\arcsin \frac{1 - x^2}{1 + x^2} \right)$

(g) $D_x \left(\arctan \frac{x+1}{x-1} + \arctan x \right)$

(c) $D_x (\arctan(x + \sqrt{1 + x^2}))$

(h) $D_x \left(\frac{\arcsin x}{\arctan x} \right)$

(d) $D_x \left(\arctan \frac{1+x}{1-x} \right)$

(i) $D_x (\arcsin(\arcsin x))$

(e) $D_x \left((\arcsin(x^2))^{-2} \right)$

7. Evaluate:

(a) $D_v \sin x$, where $v = \cos x$.

(b) $D_u \sqrt{1 - x^2}$, where $u = x^2$.

(c) $D_v (2 + 3 \cos^2 x)$, where $v = \sin x$.

8. Compute the limits of each of the following ratios.

$$(a) \lim_{x \rightarrow a} \frac{\sqrt{x^2 - 1} - \sqrt{a^2 - 1}}{x - a}$$

$$(b) \lim_{x \rightarrow a} \frac{(\arccos x)^2 - (\arccos a)^2}{x - a}$$

9. If $f(x) = (Ax + B)\sin x + (Cx + D)\cos x$, determine the value of constants A, B, C, D such that for all x , $f'(x) = x \sin x$.

10. If $g(x) = (Ax^2 + Bx + C)\sin x + (Dx^2 + Ex + F)\cos x$, determine the value of constants A, B, C, D, E, F such that for all x , $g'(x) = x^2 \cos x$.

11. Determine the derivative $g'(x)$ in terms of $f'(x)$ if:

$$(a) g(x) = f(x^3) + f(x^{-1/3})$$

$$(b) g(x) = f(\sin^2 x) + f(\cos^2 x)$$

$$(c) g(x) = f(\arcsin x) + f(\arctan x)$$

12. Prove that the derivative of an even function is odd and vice versa (it is assumed that the derivative exists).

13. Show that it is impossible to find polynomials p and q such that:

$$(a) Dp = \frac{1}{x}$$

$$(b) D \frac{1}{p} = \frac{1}{x}$$

$$\wedge (c) D \frac{p}{q} = \frac{1}{x}$$

4-7. Notation.

There are several commonly used notations for the derivative. Each of these is valuable in an appropriate context. The notation of Leibniz, in particular, will be convenient in the application of theorems on the differentiation of inverses and composite functions.

We have already used four notations for the derivative of a function at x . Consider first the three representations of the derivative,

$$(1) \quad f'(x) = \lim_{z \rightarrow x} \frac{f(z) - f(x)}{z - x} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}.$$

The notations (1) do have the virtue of complete precision: independently of context we see immediately that a specific function f is being differentiated at a specific point x . This precision was desirable for logical clarity in our development of the foundations of our subject.

In more complex situations a completely explicit notation may be a barrier to understanding rather than a help, simply because the complexity of the notation conceals our pattern of thought.

Before we examine other notations, let us review the differences in usage of the notations we already have. The two limit notations place emphasis on the numerical value of the derivative at the point x , and they are used interchangeably. The other notation places emphasis on the function

$$f' : x \rightarrow \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

The prime is usually reserved for use with the abstract designation f, g, \dots , of a function as in f', g', \dots . The symbol D_x is generally used when the function is given by an explicit name or formula, for example,

$$D_x \cot x, D_x(x^2 + ax + a^2).$$

(Here it is understood, unless the contrary is explicitly stated, that all symbols other than x appearing in the expression for the function are constants; thus,

$$D_x(x^2 + ax + a^2) = 2x + a.$$

If there is no possibility of confusion the subscript x is often omitted as in $D \cot x$.) This notation omits reference to the specific point where the derivative is being taken. The symbol D_x is an operator which when applied

to an expression giving the value $f(x)$ (for all x in a suitable domain) yields an expression giving the value $f'(x)$ of the derivative f' . Thus for

$$f : x \longrightarrow x^2$$

we have

$$D_x f(x) = D_x x^2 = 2x.$$

This must be understood as the statement that, for all x , the derivative of $f : x \longrightarrow x^2$ is $f' : x \longrightarrow 2x$. The insertion of a specific value of x in this statement makes nonsense of the initial clause, "for all x ." If we wish to make the reference to the point at which the derivative is being evaluated we must make explicit mention of it in context or invent a special notation,* for example,

$$f'(a) = D_x f(x) \Big|_{x=a}.$$

The further notations abbreviate the explicit notations by omitting the references to the function f and to the point at which the derivative is evaluated. The first of these notations parallels a common abbreviated mode of expression. We say, " y is a function of x " meaning that there exists a function

$$f : x \longrightarrow y$$

which maps each number x in a certain domain onto a value y in a certain range. This expression is appropriate when we wish to call attention to the existence of such a relation but are not impelled to name the function or define it explicitly. It is suggestive in relation to the representation of the function by its graph, the set of points (x, y) where $y = f(x)$. We use a parallel notation for the derivative,

$$f' : x \longrightarrow y'$$

and say that y' is the derivative of y , meaning that there exists a

*We cannot simply insert a value of x where x appears after the function symbol. For example, given $f : x \longrightarrow x^2$ we have $D_x f(3) = D_x 9 = 0$, the expression for the derivative of the constant function $x \longrightarrow 9$. At the same time

$$D_x f(x) \Big|_{x=3} = f'(3)$$

function f such that $y = f(x)$ and $y' = f'(x)$. Sometimes a dot is used to indicate a derivative; e.g., \dot{y} instead of y' . Clearly, to use such abbreviated notations and ways of expression we must have a context in which they are intelligible and represent a genuine convenience. We shall find them so in the next section on implicitly defined functions and their derivatives, and later in the applications.

A slightly, but significantly, more explicit notation was introduced by Leibniz. We set $y_0 = f(x_0)$ and $y = f(x)$ and write the derivative in the form

$$f'(x_0) = \lim_{x \rightarrow x_0} \frac{f(x) - f(x_0)}{x - x_0} = \lim_{x \rightarrow x_0} \frac{y - y_0}{x - x_0}.$$

Leibniz introduced the "difference" notation

$$\Delta x = x - x_0$$

$$\Delta y = y - y_0.$$

Here Δx is a single symbol with the same meaning as the symbol h used in (1), and Δy is a single symbol which represents the difference in the function values corresponding to the difference Δx between x and x_0 . In this notation, we have

$$f'(x_0) = \lim_{\Delta x \rightarrow 0} \frac{\Delta y}{\Delta x}.$$

Further, in writing the derivative, Leibniz used a parallel notation

$$\frac{dy}{dx} = \lim_{\Delta x \rightarrow 0} \frac{\Delta y}{\Delta x}.$$

Again, references to the function f and the point x_0 where the derivative is being taken are lacking and must be supplied from context. However, Leibnizian notation is slightly more explicit than Newtonian: there is not only a reference, y , to the range of f but also a reference, x , to the domain.

The symbol $\frac{dy}{dx}$ does not represent a ratio, but the limit of a ratio.

The parallelism with the notation for fractions may seem awkward at first,

* The dotted notation was introduced by Newton. It has been called "fly speck" notation by generations of irreverent teachers and students.

but it is actually singularly apt. Although the symbol $\frac{dy}{dx}$ represents a derivative and not a fraction, the rule for differentiating a composition (Theorem 4-6) permits us to handle quotients and products of these symbols formally as though they were fractions. For suitable functions

$$f : x \longrightarrow y$$

and

$$g : y \longrightarrow z$$

Theorem 4-6 states that

$$D_x gf(x) = g'(f(x)) f'(x).$$

In Leibnizian notation this becomes

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}.$$

We see, then, that the theorem has the formal appearance of a cancellation of fractions. More generally, consider a chain of compositions as indicated in Figure 4-7, where $f_1(x) = y_1$, $f_2(y_1) = y_2$, ..., $f_n(y_{n-1}) = y_n$ and $y_1 = \phi_1(x) = f_1(x)$, $y_2 = \phi_2(x) = f_2 f_1(x)$, $y_3 = \phi_3(x) = f_3 f_2 f_1(x)$,

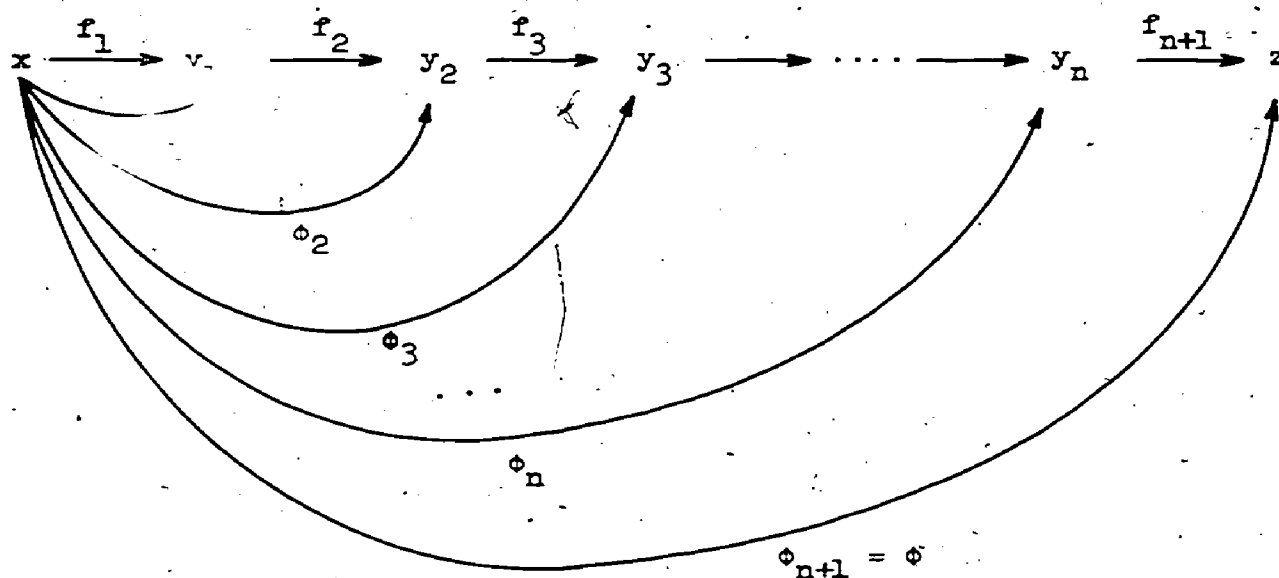


Figure 4-7

We can generalize Theorem 4-6 to this situation by applying the rule repeatedly to the compositions

$$\phi_{k+1}(x) = f_{k+1}(\phi_k(x)), \quad (k = 1, 2, \dots, n).$$

We obtain

$$\phi'_{k+1}(x) = f'_{k+1}(\phi_k(x)) \cdot \phi'_k(x).$$

Thus

$$\phi'_2(x) = f'_2(\phi_1(x)) \cdot f'_1(x),$$

$$\phi'_3(x) = f'_3(\phi_2(x)) \cdot \phi'_2(x)$$

$$= f'_3(f_2 \phi_1(x)) \cdot f'_2(\phi_1(x)) \cdot f'_1(x).$$

This generalization of Theorem 4-6 is the Chain Rule of differentiation. The proliferation of compositions is already confusing and clearly, as we go on, it will become intolerable. The mere writing of $\phi'(x)$ for large n becomes a problem. By sacrificing explicitness, Leibnizian notation resolves the problem:

$$\phi'(x) = \frac{dz}{dx} = \frac{dz}{dy_n} \frac{dy_n}{dy_{n-1}} \frac{dy_{n-1}}{dy_{n-2}} \cdots \frac{dy_2}{dy_1} \frac{dy_1}{dx}.$$

In this notation we reveal the chain of mappings,

$$x \longrightarrow y_1 \longrightarrow y_2 \longrightarrow \cdots \longrightarrow y_n \longrightarrow z$$

and at the same time exhibit the structure representing the derivative of the composition ϕ as a product of derivatives. At the same time we omit reference to the functions which define the mappings and to the points at which the several derivatives are to be taken. This information has to be supplied from context.

As a further evidence of the aptness of Leibnizian notation we observe that for inverse functions f, g , where

$$f: x \longrightarrow y$$

and

$$g: y \longrightarrow x,$$

the relation

$$g'(y) = \frac{1}{f'(x)}$$

appears in Leibnizian notation as

$$\frac{dx}{dy} = 1 / \frac{dy}{dx}.$$

It is sometimes convenient to supply the omissions of Leibnizian notation as follows: We consider $\frac{d}{dx}$ as an operator identical with D_x . We then write

$$\frac{dy}{dx} = \frac{d}{dx} f(x) = D_x f(x)$$

to supply the reference to f and write

$$\left. \frac{dy}{dx} \right|_{x=x_0} = \left. \frac{d}{dx} f(x) \right|_{x=x_0} = f'(x_0)$$

to supply the reference to the point at which the derivative is taken, e.g.,

$$\left. \frac{d}{dx} x^2 \right|_{x=3} = 6.$$

The successive higher derivatives are written $D^2 f = f''$, $D^3 f = f'''$, $D^4 f = f^{(4)}$, $D^5 f = f^{(5)}$. The Roman superscript notation becomes cumbersome for high orders and it becomes more convenient to use Hindu-Arabic numerals parenthetically as $D^{-3} f = f^{(13)}$. Thus the n -th derivative of f is written $D^n f = f^{(n)}$. It is also a useful convention to define the zero-order derivative of f as f itself, $f^{(0)} = f$.

In Leibnizian notation we write

$$D^n y = \left(\frac{d}{dx} \right)^n y = \frac{d^n y}{dx^n}.$$

Exercises 4-7

1. Let $y = \sin x$ and $x = t^2 + \frac{1}{t}$. Find $\frac{dy}{dt}\bigg|_{t=1}$ and $\frac{dy}{dx}\bigg|_{x=1}$.

2. Let $y = f(x)$ and $x = h(t)$. Express $\frac{dy}{dt}\bigg|_{t=t_0}$ in terms of t_0 .

3. Let $y = f(x)$, $x = h(t)$, $x_0 = h(t_0)$.

Using Theorem 4-6 show that

$$\frac{dy}{dx}\bigg|_{x=x_0} = \frac{\frac{dy}{dt}\bigg|_{t=t_0}}{\frac{dx}{dt}\bigg|_{t=t_0}}$$

4. Find the following:

(a) $D_x \sin x\big|_{x=0} + D_x \sin x\big|_{x=\pi/4}$

(b) $D_x(x^2 + \sin a \sin x)\big|_{x=5\pi/3}$

(c) $\frac{d}{dx}(x^2 - a^2)\big|_{x=a}$

(d) $D_x(f(a)\sin x + f(x)\sin a + f(x)\sin x)\big|_{x=a}$

5. Let $y = f(t)$, $w = g(t)$, $t = h(x)$, $z = \frac{y}{w}$.

(a) Using Leibnizian notation, find $\frac{dz}{dx}$ in terms of $\frac{dy}{dt}$, $\frac{dw}{dt}$, and $\frac{dt}{dx}$.

(b) Using (a) express $\frac{dz}{dx}\bigg|_{x=x_0}$ in terms of f' , g' , and h' .

4-8. Implicitly Defined Functions.

A function which is described in terms of rational operations on, and compositions and inverses of, known functions is said to be defined explicitly. No matter how complicated the description, if it is explicitly defined in terms of differentiable functions we know how to differentiate the function. You should, if pressed, be able to differentiate the explicit concoction.

$$(1) \quad y = \arctan \left[-1 + \sqrt{1 + \sqrt[3]{\frac{\sin^2 x}{x^2} + 1}} \right]^{1/2}$$

It often happens that a function is defined indirectly or implicitly. Thus the conditions

$$(2) \quad (\tan^4 y + 2 \tan^2 y)^3 - \frac{\sin^2 x}{x^2} = 1, \quad 0 < y < \frac{\pi}{2}$$

determine y as a function of x .

Sometimes we can find explicit expressions representing functions defined implicitly. This is the case for (2), which has the explicit solution (1). We put the implicit relationship (2) in the form

$$\tan^4 y + 2 \tan^2 y - \sqrt[3]{\frac{\sin^2 x}{x^2} + 1} = 0$$

and recognize that this is a quadratic equation for $\tan^2 y$. Solving, we obtain

$$\tan^2 y = -1 + \sqrt{1 + \sqrt[3]{\frac{\sin^2 x}{x^2} + 1}}$$

where the positive square root has been taken since $\tan^2 y$ is positive. Taking the square root and then the arctangent of both sides gives (1) since $0 < y < \frac{\pi}{2}$.

In other cases there is either no equivalent explicit definition of a function defined implicitly or it is very difficult to obtain one. An example is provided by the relation

$$(3) \quad x^2 \arctan z + z = \sin x.$$

This equation determines a unique value z for every number x ; that is, it defines a function $x \rightarrow z$ but we are unable to obtain an explicit expression for z .

It is easy to see that (3) defines z as a function of x . For any given number x , $\phi : z \rightarrow x^2 \arctan z + z$ is a continuous function and has arbitrarily large values, both positive and negative. Hence, by the Intermediate Value Theorem, there is some value z for which (3) is satisfied; since ϕ is an increasing function, that value must be unique, and the function $x \rightarrow z$ is defined implicitly by relation (3).

For the function defined by (1) we know that we can differentiate y but the execution of the differentiation would be a punishment. A much more convenient way to find the derivative is to start from (2). Applying the chain rule and other techniques of differentiation we obtain

$$3(\tan^4 y + 2 \tan^2 y)^2 (4 \tan^3 y + 4 \tan y) \sec^2 y \cdot \frac{dy}{dx} - \frac{x^2 \cdot 2 \sin x \cos x - 2x \sin^2 x}{x^4} = 0,$$

which is easily solved for $\frac{dy}{dx}$.

It is true that the formula obtained in this way will itself be somewhat implicit, since it will express $\frac{dy}{dx}$ in terms of both x and y , unlike the one we could have obtained by differentiating (1) directly, where only x would have appeared on the right side. We can still get a formula involving x alone if we want it, by using (1) to eliminate y , but it is clearly more convenient to write y instead of the complicated expression it represents. For most purposes, we do not need the completely explicit formula for the derivative. If we wish to find the value $\frac{dy}{dx}$ for a specified value of x , for instance, we can first compute the corresponding value y (explicitly from (1) in this case, but by numerical approximation in most practical problems), and then compute $\frac{dy}{dx}$ from the shorter formula.

From (3) we obtain no explicit formula for z in the first place. But we can still obtain a formula for $\frac{dz}{dx}$ by implicit differentiation. Thus, if z is a differentiable function of x , we may apply the rules of differentiation and obtain

$$2x \arctan z + x^2 \frac{1}{1+z^2} \frac{dz}{dx} + \frac{dz}{dx} = \cos x$$

or

$$(4) \quad \frac{dz}{dx} = \frac{\cos x - 2x \arctan z}{\frac{x^2}{1+z^2} + 1}.$$

If we wish to evaluate this for a specific x , we will first have to find z from (3), probably by some approximate numerical technique.

We emphasize that we have not shown that (4) holds, merely that if $\frac{dz}{dx}$ exists it must have the value given by (4). There is in fact a theorem which applies under rather general conditions (which covers the present case and most of those that arise in practice) that if an equation defining a function implicitly can be formally differentiated and the result solved for the derivative of the function, then the derivative of the function exists and has the value found. To prove, or even precisely state, this theorem would take us too far afield (Appendix 5); hereafter we shall use implicit differentiation freely to solve problems, without each time reiterating the warning that the derivative has not been proved to exist.

That we cannot solve for the derivative at every point even though the function is well defined is illustrated by the example

$$(5) \quad u^5 + x^2 u = x$$

which defines u unambiguously for each x . Implicit differentiation yields

$$(5u^4 + x^2) \frac{du}{dx} + 2xu = 1,$$

which can be solved for $\frac{du}{dx}$ everywhere except where $5u^4 + x^2$ vanishes. Since from (5) we have $u = 0$ when $x = 0$, we cannot solve for $\frac{du}{dx}$ at $x = 0$. In fact, u is not differentiable at $x = 0$.

Even if a function is differentiable at a given point the method may fail. For instance, consider the implicit definition.

$$(6) \quad v^5 + v^3 = x^3.$$

As before, at $x = 0$ we have $v = 0$ and no solution for $\frac{dv}{dx}$ from the implicitly differentiated result

$$(5v^4 + 3v^2) \frac{dv}{dx} = 3x^2.$$

In this case, however, there is a derivative at $x = 0$, and we can find it by writing (6) in the equivalent form

$$v(v^2 + 1)^{1/3} = x$$

and then differentiating:

$$[(v^2 + 1)^{1/3} + \frac{2}{3} v^2 (v^2 + 1)^{-2/3}] \frac{dv}{dx} = 1.$$

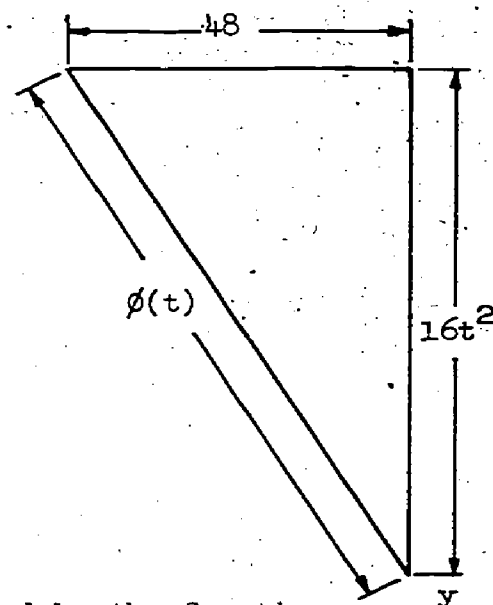
This gives $\frac{dv}{dx} = 1$ at $x = 0$.

Exercises 4-8

- For positive x , if $y = x^r$, where r is a rational number, say $r = \frac{p}{q}$ (p, q integers), then $y^q = x^p$. Assuming the existence of the derivative $D_x y$, derive the formula $D_x y = rx^{r-1}$ using implicit differentiation and the differentiation formula $D_x x^n = nx^{n-1}$, for integral n .
- For each of the following, find $D_x y$ without solving for y as a function of x .
 - $5x^2 + y^2 = 12$
 - $2x^2 - y^2 + x - 4 = 0$
 - $y^2 - 3x^2 + 6y = 12$
 - $x^3 + y^3 - 2xy = 0$
- For each of the following use implicit differentiation to find $D_x y$:
 - $x^2 = \frac{y-x}{y+x}$
 - $x^2 y + xy^2 = x^3$
 - $x^m y^n = 10$ (m, n integers)
 - $\sqrt{xy} + x = y^{-1}$
- Use implicit differentiation to find $D_y x$.
 - $x\sqrt{y} + y\sqrt{x} = a\sqrt{a}$
 - $2x^2 + 3xy + y^2 + x - 2y + 1 = 0$
 - $(x+y)^{1/2} + (x-y)^{1/2} = 4$
 - $3x^2 + x^2 y^2 = y^4 + 5$
 - $4x^2 + 3xy - 7y^2 = 0$
- For each equation, find the slope of the curve represented, at the stated point.
 - $2x^2 + 3xy + y^2 + x - 2y + 1 = 0$ at the point $(-2, 1)$.
 - $x^3 + y^2 x^2 + y^3 - 1 = 0$ at the point $(1, -1)$.
 - $x^2 - x\sqrt{xy} - 6y^2 = 2$ at the point $(4, 1)$.
 - $x \cos y = 3x^2 - 5$ at the point $(\sqrt{2}, \frac{\pi}{4})$.

6. For each equation, find the slope of the curve represented at the point or points where $x = y$. Give a geometric explanation for these results.
- $x^3 - 3axy + y^3 = 0$
 - $x^m + y^m = 2$
 - $x^2 + y^2 = 2axy$
7. Find $D_x y$ by implicit differentiation.
- $a \sin y + b \cos x = 0$
 - $x \cos y + y \sin x = 0$
 - $\sin xy = \sin x + \sin y$
 - $\csc(x + y) = y$
 - $x \tan y - y \tan x = 1$
 - $y \sin x = x \tan y$
 - $xy + \sin y = 5$
8. If $0 < x < a$, then the equation $x^{1/2} + y^{1/2} = a^{1/2}$ defines y as a function of x . Assuming the existence of the derivative, show without solving for y that $f'(x)$ is always negative.
9. Assuming that $D_x y = D_y x = 0$ (i.e., x and y are independent), find the following.
- $D_x(x^2 + xy + \cos y)$
 - $D_x(y^2) + D_y(x^2)$
 - $D_x(x^2) + D_y(y^2)$
 - $D_x f(xy) + D_y f(x)$
 - $D_x(xy)^2$
10. Let c_1 and c_2 be two curves which intersect at the point (x_0, y_0) and let the slopes of c_1 and c_2 at (x_0, y_0) be m_1 and m_2 , respectively. If the product $m_1 m_2$ equals -1 , we may say that the curves c_1 and c_2 are orthogonal.
- Show that the lines with equations $4y - 3x - 40 = 0$ and $3y + 4x + 15 = 0$ are orthogonal.
 - Show that the circle $x^2 + y^2 = r^2$, r constant, is orthogonal to the line $y = mx$, m constant.

11. Find the points of intersection of the ellipse $x^2 + 10y^2 = 10$ and the hyperbola $x^2 - 8y^2 = 8$, and the slopes of the curves at these points of intersection. Show that the curves are orthogonal.
12. Show that the family of curves $y^2 = 4a(x + a)$ is self-orthogonal, i.e., each two members of the family $y^2 = 4a_1(x + a_1)$ and $y^2 = 4a_2(x + a_2)$ that intersect, necessarily intersect at right angles.
- A13. For what values of k will there be exactly one line passing through the point $(0, k)$ and orthogonal to the parabola $y = x^2$? For what values of k will there be exactly three orthogonal lines?
14. A ball dropped out of a window falls $16t^2$ feet in t seconds. An observer is watching from another window at the same height 48 feet away. At what rate is the distance of the ball from the observer increasing two seconds after the ball is dropped?
- (a) Write an equation which implicitly defines the distance $y = \phi(t)$ between the observer and the ball at time t .
- (b) Use implicit differentiation to answer the question of the problem.
15. (a) Given that simple harmonic motion is described by the function $\rho : t \rightarrow \sin(\omega t + c)$ where ω and c are constants. Find the velocity at time $t = t_0$.
- (b) Simple harmonic motion may also be described by the function $\rho : t \rightarrow \cos(\omega t + c)$ where ω and c are constants. Find the velocity at time $t = t_0$.
- (c) In what sense are the motions in (a) and (b) the same?
16. If a simple harmonic motion is described by the function $\rho : t \rightarrow A \sin \omega t + B \cos \omega t$ where A , B , and ω are constants, determine the maximum velocity.



Miscellaneous Exercises

1. Evaluate:

(a) $D(x + \frac{1}{x})^{1/2}$.

(k) $D_v(\sin x \cdot \cot x)$, where $v = \cos x$.

(b) $D((\arcsin \sqrt{x})^2)$.

(l) $D(\sin(x^{-1/2}) - \cos(x^{1/2}))$.

(c) $D\sqrt{3x^2 - 1} \Big|_{x=3/2}$.

(m) $D_v(\frac{1 - \sin x}{1 + \cos x})$, where $v = \cos \frac{x}{2}$.

(d) $D_u \sqrt{x - x^2}$, where $u = x^2$.

(n) $D(\frac{2+t}{3t^2 - 1})$.

(e) $D\sqrt{x + \sqrt{x}}$.

(o) $D(\frac{(2x-3)^2}{x^3 - 1})$.

(f) $D(x^3 - 2)^6$.

(p) $D(2x^3 + 5x^2 - x + 2)^{10}$.

(g) $D\sqrt[3]{x^2 - 3}$.

(q) $D_v(\frac{1-x}{x^2})$, where $v = \arcsin x$.

(h) $D_u(x^2 - x^{-1/2} + x^{-2})$,

(r) $D \cos \sqrt{1 - x^2}$.

where $u = \sqrt{x}$.

(i) $D(\frac{x}{x + \sqrt{a^2 - x^2}})$,

(s) $D(\arcsin(x^2) + \arccos(\frac{\pi}{2} - x^2))$.

a constant.

(j) $D\sqrt{x+1}$.

(t) $D(\tan(x^{-1} - x))$.

2. In Section 2-4 we defined the velocity of an object whose location on a straight line at time $t = t_0$ is given by $s = \phi(t)$ as the limit of the ratio

$$\frac{\phi(t) - \phi(t_0)}{t - t_0},$$

and in Section 2-5 we observed that this limit is the value of the derivative ϕ' at $t = t_0$. Experimentally it has been established that the distance covered in time t by a freely falling body is proportional to t^2 , and therefore it can be represented by the function $\phi: t \rightarrow ct^2$, where c is a positive constant. Show that the velocity of a freely falling body is directly proportional to the time.

3. Suppose a projectile is ejected at a point P which is 20 feet above the ground with initial velocity of v_0 feet per second. Neglect friction and assume that the projectile moves up and down in a straight line. Let $\theta(t)$ denote the height (above P) in feet that the projectile attains t seconds after ejection. Note that if gravitational attraction were not acting on the projectile, it would continue to move upward with a constant velocity, traveling a distance of v_0 feet per second, so that its height at time t would be given by $\theta(t) = v_0 t$. We know that the force of gravity acting on the projectile causes it to slow down until its velocity is zero and then travel back to the earth. On the basis of physical experiments the formula $\theta(t) = v_0 t - \frac{g}{2} t^2$, where g represents the force of gravity, is used to represent the height (above P) of the projectile as long as it is aloft. Note that $\theta(t) = 0$ when $t = 0$ and when $t = \frac{2v_0}{g}$. This means that the projectile returns to the initial 20 foot level after $\frac{2v_0}{g}$ seconds.

- (a) Find the velocity of the projectile at $t = t_0$ (in terms of v_0 and g).
- (b) Sketch the s vs. t and the v vs. t curves on the same set of axes.
- (c) Compute (in terms of v_0) the time required for the velocity to drop to zero.
- (d) What is the velocity on return to the initial 20 foot level?
- (e) Assume that the projectile returns to earth at a point 30 feet below the initial take-off point P . What is the velocity at impact?

Chapter 5

APPLICATIONS OF THE DERIVATIVE

5-1. Introduction.

In the preceding chapter we were concerned primarily with the concept of derivative at a point. The derivative $f'(a)$ of a function f at a point a is an example of a local property of f : it is defined as

$$f'(a) = \lim_{h \rightarrow 0} \frac{f(a+h) - f(a)}{h},$$

so that the derivative at a is defined by the values of f in any neighborhood of a , no matter how small. By contrast, in the applications, we are often concerned with global properties of f ; for example, whether f is increasing throughout its domain, or whether M is the largest value $f(x)$ for all x in its domain. A global property of f involves the values of f in its entire domain, a local property involves only values in the neighborhood of some given point. For example, the property that $f(x)$ is bounded below on $[a,b]$ by a given constant is global. The property $f'(x) \geq 0$ at a given point is local.

One of the problems of this chapter is to relate global properties of f to local conditions satisfied by f at each point of its domain. For example, if f is continuous on $[a,b]$ and differentiable on (a,b) , the global property $f(x) \geq f(a)$ for all x in $[a,b]$ is established if $f'(x) \geq 0$ for all x in (a,b) .

In general, we explore the properties of functions differentiable on an interval. We focus our attention upon the function

$$Df = f'$$

rather than upon the values of the derivative at individual points. We also consider the functions obtained by iterating the operation of differentiation; in particular, we shall investigate properties related to the second derivative,

$$D(Df) = D^2f = f'',$$

and shall also have occasional use in this chapter for derivatives of higher order.

Our major aim in this chapter is to see how the derivative f' gives us information about the function f and its graph. We already have obtained one such item of information in Theorem 3-6d: because f must be defined and continuous where it is differentiable, each point x on the domain of f' is also a point of continuity of f . From this fact, upon any closed interval in the domain of f' , we obtain the global properties of f given by the extreme and intermediate value theorems of Section 3-7. Further, we shall next see how knowledge of the derivative of f on an interval enables us to locate the extreme values in that interval. In Section 5-5 we shall see how information about the second derivative f'' , the derivative of f' , gives information about f' and also about the function f itself.

Exercises 5-1

1. (a) Obtain an expression for the first and higher derivatives of x^n where n is a natural number.
- (b) Do the same for $x^{p/q}$ where $\frac{p}{q}$ is rational, p and q relatively prime. What is the domain of f' when $f(x) = x^{1/2}$? For what values $\frac{p}{q}$ is the domain of f' different from that of f ? Answer the same question for higher derivatives of $x^{p/q}$.

2. What is the twenty-third derivative of

- (a) $\sin x$? (c) $\sin 2x$?
 (b) $\cos x$? (d) $5x^{17} - \sqrt{2}x^{10} - 83x^3 - 70425$?

3. Find the n -th derivatives of the following functions:

(a) $f : x \rightarrow (ax + b)^{-1}$.

(b) $f : x \rightarrow \sin x$.

(c) $f : x \rightarrow \cos(ax + b)$.

(d) $f : x \rightarrow \sin^2 x$.

(e) $f : x \rightarrow \cos^3 2x$.

(f) $f : x \rightarrow x^n - [x]^n$.

(g) $f : x \rightarrow \frac{1}{x^2 - a^2}$.

(h) $f : x \rightarrow \frac{1}{x^2 + 1}$.

4. Let f be a function defined for all real values of x with the property

$$f(a + b) = f(a) \cdot f(b), \text{ for all real numbers } a \text{ and } b.$$

- (a) Show that either $f(0) = 1$ or $f(x) = 0$ for all x .
- (b) If $f(0) = 1$, show that $f(x) \neq 0$ for all x .
- (c) If, additionally, $f(0) = 1$ and f has a derivative at $x = 0$, show that f' exists for all x and that $f'(x) = f'(0) \cdot f(x)$.

5. If $F(x + y) = F(x) + F(y)$, where F is continuous at one point and defined for all x , show that $F(x)$ is everywhere differentiable.
6. Given that the function F is defined for all x , that $F'(x)$ exists for $0 \leq x < a$, and that $F(x + a) = kF(x)$ for all x (a, k constant). Show that F is everywhere differentiable and sketch some possible graphs of F assuming
- (a) $F(0) = 0$, $F'(0) = 0$.
 - (b) $F(0) = 0$, $F'(0) = 1$.
 - (c) If $F(0) = 0$ and the graph of F is tangent to the y -axis show that F is differentiable for all $x \neq 0$ and sketch some possible graphs of F .

5-2. The Derivative at an Extremum.

In the problem of Section 1-1, we attempted to determine the dimensions of a box of greatest volume V that satisfies certain postal regulations: we sought to maximize the function $x \rightarrow 72x^2 - 4x^3$. Another problem of similar mathematical structure might be to determine the length L of the shortest path along which two points can be connected so as to meet a given curve somewhere on the path (the reflection problem of a later chapter on geometrical optics). Here we seek the minimum of the function L . A value which is either a maximum or a minimum for a function f is called an extremum of f . *

(i) Location of an extremum on a closed interval.

We consider a function f differentiable on the open interval $a < x < b$ and continuous on the closed interval $a \leq x \leq b$. Since f is continuous on the closed interval it assumes a maximum value M there (Theorem 3-7b). That is, for all values x such that $a \leq x \leq b$ we have $f(x) \leq M$ and for at least one value u we have $f(u) = M$. The possibility that u is an end-point of the interval must always be considered. If u is an interior point of the interval then, as we indicated in Section 1-1, for $f(u)$ to be a maximum, $f'(u)$ must be 0. Now we prove that fact.

THEOREM 5-2a. Let $f(u)$ be an extremum of f for some value u in the interior of the domain of f . If $f'(u)$ exists, then $f'(u) = 0$.

Proof. Suppose $f'(u) \neq 0$. Then either $f'(u) > 0$ or $f'(u) < 0$. Suppose $f'(u) > 0$; then, since $f'(u) = \lim_{x \rightarrow u} \frac{f(x) - f(u)}{x - u}$, by Lemma 3-4, there is a neighborhood I of u such that

$$r(x) = \frac{f(x) - f(u)}{x - u} > 0, \quad \text{for } x \text{ in } I.$$

Since u is an interior point of $[a, b]$ there will then be points α, β in the interval for which $\alpha < u < \beta$ and both $r(\alpha)$ and $r(\beta)$ are positive. Consequently,

$$f(\alpha) < f(u) < f(\beta),$$

* Maximum, minimum, and extremum (-um becomes -a in plural) are words of Latin origin meaning greatest, smallest, and outermost respectively.

and $f(u)$ cannot be an extremum.

There is a parallel proof for the case $f'(u) < 0$.

Corollary. If f is differentiable on the open interval and continuous on the closed interval $[a, b]$ then the extreme values are taken on at endpoints of the interval or at interior points where the derivative is zero.

Example 5-2a. To determine extreme values of $f(x) = 1 + x - x^2 - x^3$ on the interval $[0, 1]$, we obtain

$$f'(x) = (1 - 2x - 3x^2) = (1 - 3x)(1 + x).$$

We conclude that

$$f'(\frac{1}{3}) = f'(-1) = 0.$$

However, $x = -1$ is not a point of the interval and we conclude that the extreme values can be taken on only at the interior point $\frac{1}{3}$ or the endpoints 0 and 1. We have $f(0) = 1$, $f(\frac{1}{3}) = \frac{32}{27}$, $f(1) = 0$. We conclude that the maximum value is taken on at the interior point $\frac{1}{3}$ and the minimum value at the endpoint 1.

Example 5-2b. In the parcel post problem of Section 1-1 we dealt with the function $f : x \rightarrow 72x^2 - x^3$ on the interval $[0, \frac{72}{5}]$. The derivative $f'(x) = 144x - 12x^2$ is zero at $x = 0$ and $x = 12$. We have $f(0) = 0$, $f(12) = 3456$, $f(\frac{72}{5}) = (\frac{72}{5})^3 = 2985 \frac{123}{125}$. With this observation we have completely proved that the solution given in Section 1-1 does provide the carton of greatest capacity. We see also that the minimum value of f occurs at the endpoint $x = 0$; that the derivative happens to be zero there is irrelevant.

Example 5-2c. Consider the function $f : x \rightarrow |x|$ on the interval $[-1, 1]$. We know that $|x| \geq 0$ and that $|0| = 0$. It follows that the minimum of $|x|$ on $[-1, 1]$ is zero and is taken on at $x = 0$. Here is a case where an extreme value is taken on at an interior point, but the conditions of the corollary do not apply; $|x|$ does not have a derivative at $x = 0$. We have had to appeal to other evidence to locate the minimum. However, $|x|$ is differentiable on the open intervals $(-1, 0)$ and $(0, 1)$.

Since the derivative of $|x|$ is nowhere zero we conclude that any extreme value of $|x|$ on $[-1,1]$ must be located among the points $x = -1$, $x = 0$, $x = 1$.

With a solid theoretical foundation which was lacking in Chapter 1 we now have an effective method of solution for a broad class of extreme value problems. In summary, for a function f continuous on a closed interval we know that the extrema exist and we know that if an extremum exists at an interior point u of the interval where f is differentiable then $f'(u) = 0$. To locate the extrema of f on the interval, then, we need only try endpoints of the interval, points where f has no derivative, and points where the derivative is zero. (Most of the functions considered here are differentiable everywhere, although in some cases there may be exceptional points where the derivative does not exist.) To determine which of these points yield the extrema we may calculate the values of the function at each point of this restricted class; the largest such value is the maximum value of the function on the interval, the smallest is the minimum.

Exercises 5-2a

1. Complete the proof of Theorem 5-2a by proving that $f'(u) \neq 0$.
2. Make a careful sketch on the interval $[0,1]$ of the graph of the function $f : x \mapsto 1 + x - x^2 - x^3$ given in Example 5-2a. Does the graph confirm the conclusions of the text?
3. Determine the extreme values of the function $f : x \mapsto 44 + 4x - 13x^2 + 18x^3 - 9x^4$ and make a careful sketch of the graph of f . Compare your results with your answer to Exercise 1-1, No. 12.
4. Locate and characterize the extreme values of each of the following functions on the interval $[-1,1]$.
 - (a) $f : x \mapsto x^{2/3}$
 - (b) $f : x \mapsto |x|^{3/2}$
5. (a) Three men live on the same straight road. Where on the road should they agree to meet so that the ~~sum~~ of the distances they travel along the road from their homes to their meeting place is to be a minimum?
- (b) What is the answer if the number of men is four?
- (c) Answer the question for n men where n is any positive integer.

6. A stone wall 100 yards long stands on a ranch. Part or all of it is to be used in forming a rectangular corral, using an additional 260 yards of fencing for the other three sides. Find the maximum area which can be so enclosed.
7. A metal can with square base and open top is required to contain a gallon (231 cubic inches) of gasoline. Neglecting the thickness of the metal and the waste material in construction, find the dimensions that require the least material.
8. A right triangle with hypotenuse k is rotated about one of its legs. Find the maximum volume of the right circular cone produced.
9. Determine the lengths of the sides of a triangle of maximum area with base b and perimeter p . (Hint: Use Heron's formula for the area of a triangle: $A = \sqrt{s(s-a)(s-b)(s-c)}$ where a, b, c are the lengths of the sides, and $s = \frac{1}{2}(a+b+c) = \frac{p}{2}$.)

(ii) Local extrema.

In the preceding examples we showed how to locate the global extrema on a closed interval of a differentiable function f . Clearly an extremum of the function at an interior point of a given interval may fail to be an extremum in the interior of a larger interval. However, if $f(u)$ is an extreme value of f in some neighborhood of u , then it is an extremum in any smaller neighborhood; in this case, we call $f(u)$ a local extremum of f (some texts use the term relative extremum). From Theorem 5-2a if $f(u)$ is a local extremum then $f'(u) = 0$ provided $f'(x)$ exists. In seeking the global extrema by the method described above we shall also find the local extrema of f .

Example 5-2d. Consider the function $f(x) = 3x^4 + 4x^3 - 12x^2 + 5$. We have $f'(x) = 12x^3 + 12x^2 - 24x = 12x(x-1)(x+2)$, whence

$$f'(0) = f'(1) = f'(-2) = 0.$$

We tabulate the values of f at the zeros of the derivative and at the end-points of an interval including the zeros.

x	-3	-2	0	1	2
$f(x)$	32	-27	5	0	37

From the table we conclude that $f(-2)$ is the minimum of $f(x)$ on the interval $[-3, 2]$, $f(0)$ a maximum on $[-2, 1]$, and $f(1)$ a minimum on $[0, 2]$. Since there are no other local maxima or minima we expect the graph of f to have the appearance of Figure 5-2a.

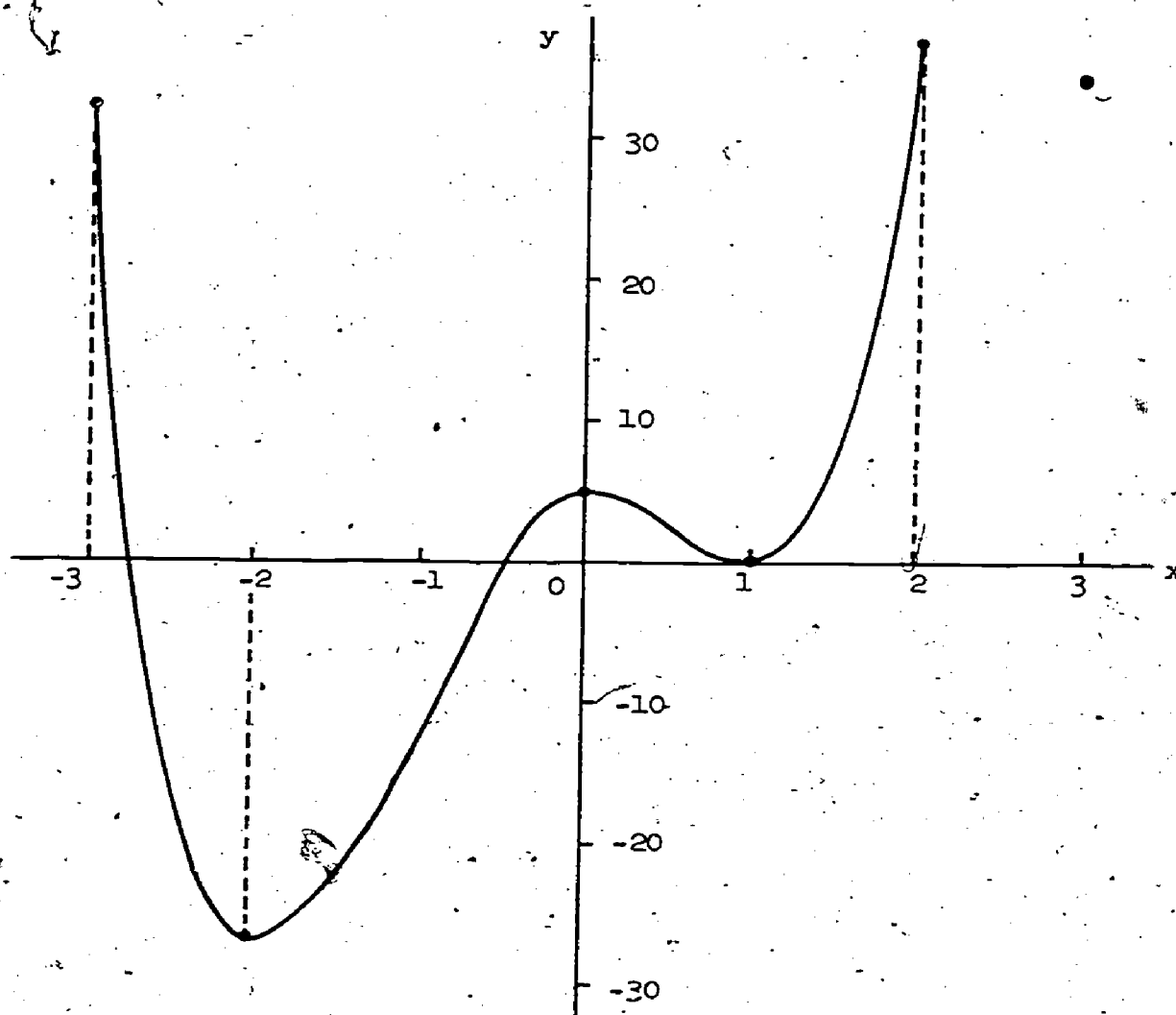


Figure 5-2a

Thus, we expect f to be a decreasing function for $x \leq -2$; $f(-2)$ is a local minimum. For $-2 \leq x \leq 0$ the function should be increasing and $f(0)$ is a local maximum. Over the interval $0 \leq x \leq 1$ the function should decrease again to the local minimum $f(1)$; for $x > 1$ the function should be increasing.

We shall prove that the inferences concerning the function of the example are correct. By determining the positions of the local extrema of the function we have been able to say a great deal about the general character of the function. We have utilized intuitively the idea that throughout the interval between successive extrema, the function must either be increasing or decreasing. We shall prove this result as Theorem 5-2b. In the course of the proof we shall utilize a result we shall need again.

Lemma 5-2. Let f be a continuous function on the closed interval $[a, b]$. If $f(a) = f(b)$, then there is a local extremum for $f(x)$ on the open interval (a, b) .

Proof. Assume no extremum exists on the interior of the interval. Since extrema exist on the closed interval by the extreme value theorem, the extrema of f must occur at the endpoints. Since $f(a) = f(b)$ the maximum and minimum of f must then be the same. It follows that f is constant on $[a, b]$. Hence $f(x)$ is an extreme value for all x satisfying $a \leq x \leq b$. This contradicts the assumption that $f(x)$ has no extremum on (a, b) .

The two simple results, Theorem 5-2a and Lemma 5-2, constitute a basis for the logical development of the rest of the chapter.

THEOREM 5-2b. If f is continuous on the closed interval $[a, b]$ and if it has no local extrema on the open interval (a, b) then f is strongly monotone* on $[a, b]$.

Proof. By Lemma 5-2, we have $f(a) \neq f(b)$, hence either $f(a) < f(b)$ or $f(a) > f(b)$. If $f(a) < f(b)$ we shall prove f is increasing on $[a, b]$. (A parallel argument would prove f is decreasing if $f(a) > f(b)$.)

Let u, v be any two points of $[a, b]$ with $u < v$. We want to show that $f(u) < f(v)$. The case $f(u) = f(v)$ is impossible in view of Lemma 5-2 and thus it will be sufficient to show that the assumption $f(u) > f(v)$ results in a contradiction.

We will first establish a simple consequence of the Intermediate Value Theorem: if α, β, γ are three points in $[a, b]$ such that

* See Definition A2-4b.

$$(1) \quad \alpha < \beta < \gamma$$

and

$$(2) \quad f(\beta) > f(\alpha), f(\beta) > f(\gamma),$$

then there exist distinct points c_1 and c_2 such that $\alpha \leq c_1 < \beta < c_2 \leq \gamma$ and $f(c_1) = f(c_2)$. (See Figure 5-2b.)

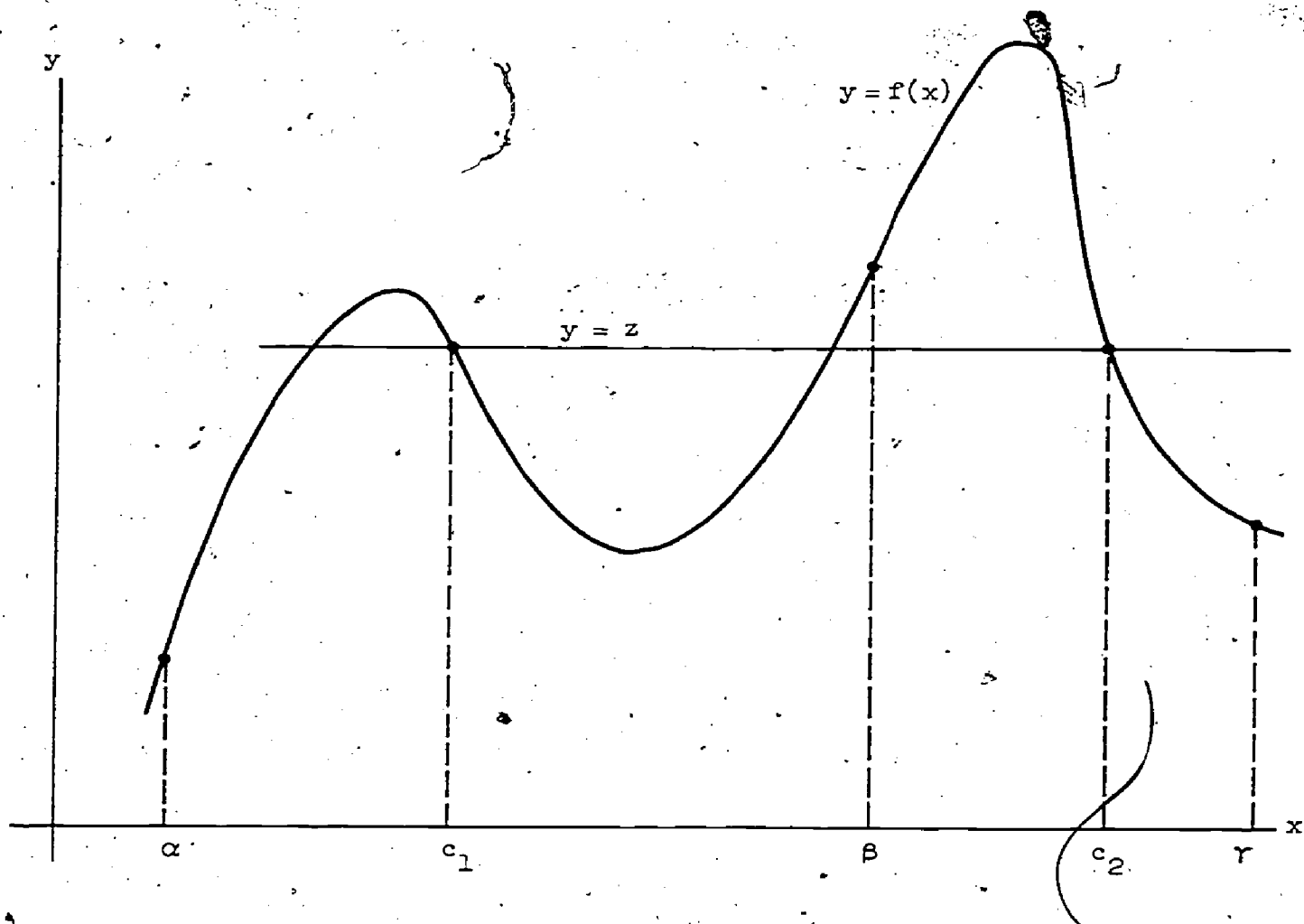


Figure 5-2b

For, let z be any value less than $f(\beta)$ but greater than both $f(\alpha)$ and $f(\gamma)$; for example, $z = \frac{1}{2} [f(\beta) + \max\{f(\alpha), f(\gamma)\}]$. Then $f(\alpha) < z < f(\beta)$, and by the Intermediate Value Theorem there exists a $c_1 \in [\alpha, \beta]$ such that $f(c_1) = z$; likewise, since $f(\gamma) < z < f(\beta)$, there exists a $c_2 \in [\beta, \gamma]$ such that $f(c_2) = z$. Finally, $c_1 \neq \beta$, $c_2 \neq \beta$, since $f(c_1) = f(c_2) = z \neq f(\beta)$.

A similar argument produces the same result in the case when the inequalities in (2) are reversed.

Next, under the assumption that $f(u) > f(v)$, we will exhibit two points c_1 and c_2 in $[a, b]$ such that $f(c_1) = f(c_2)$. It then follows, by Lemma 5-2, that f has a local extremum on (a, b) in contradiction to the hypothesis of the theorem.

Obviously, either $f(u) \leq f(a)$ or $f(u) > f(a)$. Suppose first that $f(u) \leq f(a)$. Thus $f(v) < f(u) \leq f(a) < f(b)$.

We conclude then that $a < v < b$ while $f(v) < f(a)$ and $f(v) < f(b)$. We can now apply the previous observation (with $\alpha = a$, $\beta = v$, and $\gamma = b$ and with inequalities in (2) reversed) to exhibit the desired points c_1 and c_2 .

Finally, suppose that $f(u) > f(a)$. Then $u \neq a$ so that $a < u < v$, while $f(u) > f(a)$ and $f(u) > f(v)$. Again we can apply the previous observation (with $\alpha = a$, $\beta = u$, and $\gamma = v$) to exhibit the desired points c_1 and c_2 . This concludes the proof of the theorem.

Our method for finding the global extrema of a differentiable function has a useful by-product: the local extrema also are determined. Theorem 5-2b justifies the description of the gross properties of the function given in Example 5-2d.

We know that on an open interval all the local extrema of a differentiable function f are selected by the condition $f'(x) = 0$. There may be points for which $f'(x) = 0$, however, which do not correspond to local extrema of f .

Example 5-2e. Consider $f(x) = x^3$ on the interval $[-1, 1]$. We have $f'(x) = 3x^2$, hence $f'(0) = 0$ and the derivative vanishes nowhere else. We have $f(-1) = -1$, $f(0) = 0$, and $f(1) = 1$. It is easy to prove that x^3 is an increasing function, hence that $f(0)$ is not an extremum. We conclude only that the graph of f is horizontal at $x = 0$, but this information is also useful in sketching the graph. (See Figure 5-2c.)

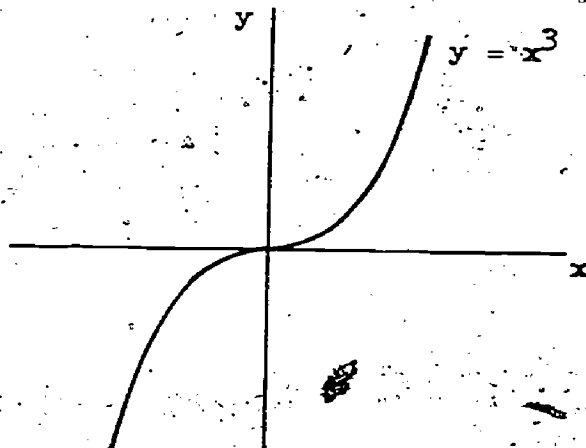


Figure 5-2c

The result of the preceding example is a particular instance of a corollary of Theorem 5-2b.

Corollary 1. Let f be continuous on the closed interval $[a,b]$ and differentiable on the open interval (a,b) . If there exists only one point u in (a,b) where $f'(u) = 0$ and if either $f(a) < f(u) < f(b)$ or $f(a) > f(u) > f(b)$, then $f(u)$ is not a local extremum.

The proof of this corollary is left to you as an exercise (Exercises 5-2b, No. 1).

To supplement Corollary 1 we give a further result.

Corollary 2. Let f be continuous on the closed interval $[a,b]$ and differentiable on the open interval (a,b) . Let there be only one point u in the open interval where $f'(u) = 0$. If $f(u) > f(a)$ and $f(u) > f(b)$ then $f(u)$ is the maximum of f on $[a,b]$. If $f(u) < f(a)$ and $f(u) < f(b)$ then $f(u)$ is the minimum of f on $[a,b]$.

The proof of Corollary 2 is left as an exercise (Exercises 5-2b, No. 1).

Example 5-2f. We apply the knowledge we have gained to find the local and global extrema of the function

$$f : x \mapsto 4x^5 - 5x^4 - 40x^3 + 100$$

on the interval $-3 \leq x \leq 4$. We differentiate and obtain

$$f'(x) = 20x^4 - 20x^3 - 120x^2 = 20x^2(x+2)(x-3)$$

Computing the values of f at the zeros of f' and at the endpoints we obtain the following table.

x	-3	-2	0	3	4
$f(x)$	-197	212	100	-413	356

Considering triples of consecutive values of f in this table in the light of Corollaries 1 and 2 to Theorem 5-2b we find that the function f increases from a local minimum at $x = -3$ to a local maximum at $x = -2$, then decreases to its global minimum at $x = 3$ and increases to its global maximum at $x = 4$. (If we were to consider the entire real axis as the domain of f then, since f' has no zeros outside the interval $(-3, 4)$, we would conclude that f is increasing for $x < -3$ and increasing for $x > 4$.) We can utilize the information of the table and a few additional plotted points to obtain an excellent idea of the behavior of the graph of f on the given interval (Figure 5-2d).

In summary, to locate the extrema of a continuous function on a closed interval, in the light of the preceding discussion we restrict our search to the endpoints of the interval, interior points where the derivative does not exist, and points where the derivative is zero.

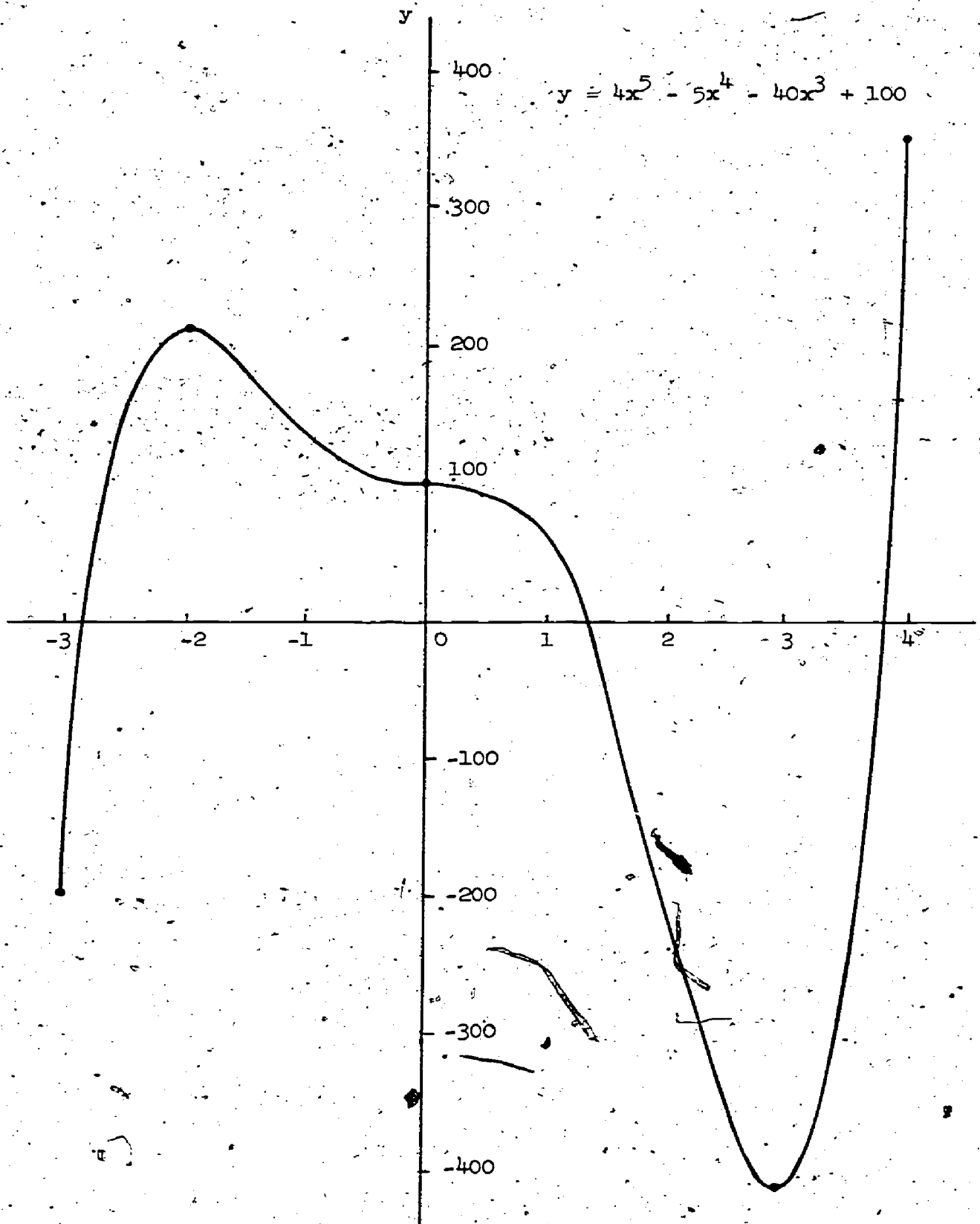



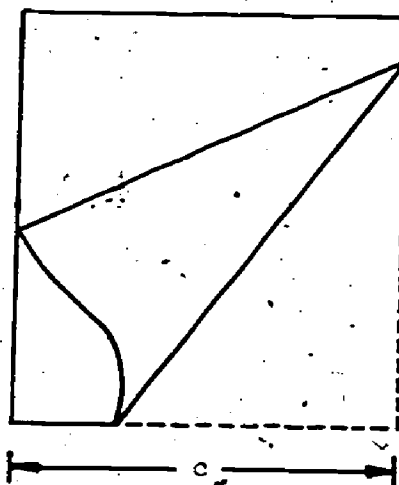
Figure 5-2d

Exercises 5-2b

1. Prove the corollaries to Theorem 5-2b.
2. For each of the following functions locate and characterize all extrema. On what intervals is the function increasing? decreasing?
 - (a) $f: x \rightarrow 4x^4 - 8x^2 + 1$
 - (b) $f: x \rightarrow x^4 - 4x^3$
 - (c) $f: x \rightarrow \frac{x^3}{1 + x^2}$
 - (d) $f: x \rightarrow \frac{x}{x^2 - 1}$
 - (e) $f: x \rightarrow \frac{x}{1 + x^2}$
3. A rectangle is inscribed in a circle of radius R . Find the rectangle of maximum area; of maximum perimeter.
4. The area of the printed text on a page is A square centimeters. The left and right margins are each c centimeters wide, and the upper and lower margins are each d centimeters. What are the most economical dimensions of the pages if only the amount of paper matters?
5. A rectangle has two of its vertices on the x -axis and the other two above the axis on the parabola $y = 6 - x^2$. What are the dimensions of such a rectangle if its area is to be a maximum?
6. A rectangular sheet of galvanized metal is bent to form the sides and bottom of a trough so that the cross section has this shape: 

If the metal is 14 inches wide, how deep must the trough be to carry the most water?
7. Find the right circular cylinder of greatest volume that can be inscribed in a right circular cone of radius r and height h .

8. The lower right-hand corner of a page is folded over so as to reach the left edge in such a way that one endpoint of the crease is on the right-hand edge of the page and the other endpoint is on the bottom edge of the page as in the figure. If the width of the page is c inches, find the minimum length of the crease.



9. What is the smallest positive value of t such that the slope of $y = 2 \sin(\frac{t}{2} - \frac{\pi}{3})$ is zero?
10. A wall h feet high stands d feet away from a tall building. A ladder L feet long reaches from the ground outside the wall to the building. Let ϕ be the angle between the ladder and the building.
- (a) Show that if the ladder touches the top of the wall,
 $L = d \csc \phi + h \sec \phi$.
- (b) Find the shortest ladder that will reach the building if $h = 8$ and $d = 24$.
11. In an experiment repeated n -times, one obtains the numbers a_1, a_2, \dots, a_n for a certain physical quantity x . What value of x should we take if we want to:
- (a) minimize the sum of the squares of the deviations, i.e.,
 $(x - a_1)^2 + (x - a_2)^2 + \dots + (x - a_n)^2$.
- (b) minimize the sum of the absolute values of the deviations, i.e.,
 $|x - a_1| + |x - a_2| + \dots + |x - a_n|$.
12. Find the maximum of $x^m y^n$ (m, n rational and > 0) if $x + y = c$ (c constant) and $x \geq 0, y \geq 0$.
13. Find the minimum of $x + y$ if $x^m y^n = k$ (k constant) and (m, n rational and > 0).

5-3. The Law of the Mean.

Until now we have used derivative of a function only to locate extrema. As we shall see in Section 5-4, however, the derivative f' in one sense determines the function f almost completely. In order to make use of this fact we need an efficient way of arguing from properties of the derivative to properties of the function. The Law of the Mean provides such a way.

(i) Statement and proof of the Law of the Mean.

In geometrical terms, the Law of the Mean states that on the arc between any two points of the graph of a differentiable function there exists a point where the curve has the same slope as the chord.* Thus, let $(p, f(p))$ and $(q, f(q))$ be any two points on the graph of a differentiable function f with $p < q$, say (see Figure 5-3a).

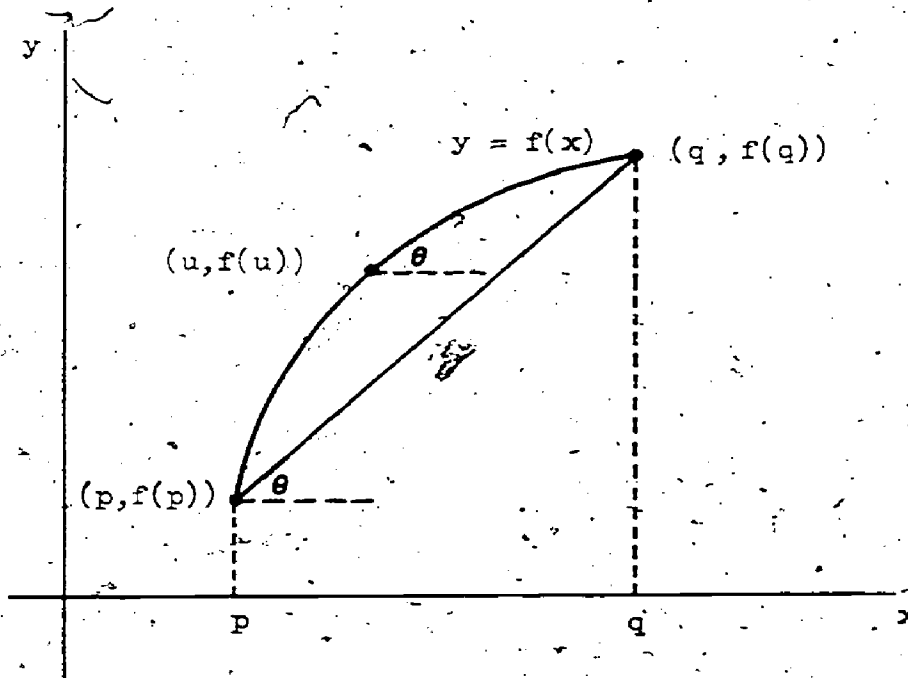


Figure 5-3a

According to the Law of the Mean there exists a point u between p and q where

*The word "mean" here signifies "average". The slope of the chord is interpreted as average rise in function value per rise in value of x . The Law of the Mean states that this average is equal to a value of the derivative at some point of the interval.

$$f'(u) = \frac{f(q) - f(p)}{q - p}$$

We can make the Law of the Mean plausible by an argument similar to that by which we found that the slope of a graph at an interior extremum is zero (Section 1-1). Take a parallel to the chord at a point $(u, f(u))$ which lies on the arc at maximum distance from the chord. Since no point of the arc lies at a greater distance from the chord, the arc cannot cross the parallel. The arc cannot meet the parallel at an angle for then it would cross; therefore the two must have the same direction at $(u, f(u))$. (See Figure 5-3b.)

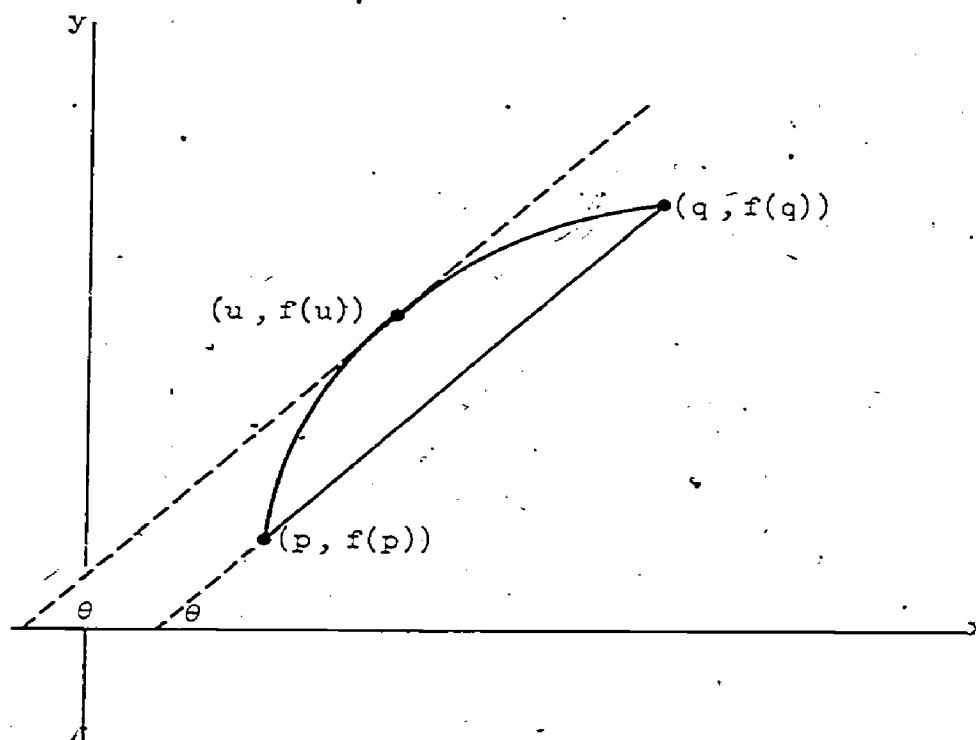


Figure 5-3b

In order to derive the Law of the Mean we first prove it for the special case in which the chord is horizontal.

Lemma 5-3. (Rolle's Theorem). If f is continuous on the closed interval $[p, q]$ and differentiable on the open interval (p, q) , and if $f(p) = f(q)$, then there is at least one point u in the open interval where $f'(u) = 0$.

Proof. From the assumption of continuity alone we have already shown that f has an extremum on the open interval (Lemma 5-2). If u is a point of (p, q) for which $f(u)$ is an extremum then since f is differentiable on the open interval we know by Theorem 5-2a that $f'(u) = 0$.

Before proving the Law of the Mean, let us examine some of the other consequences of Rolle's Theorem (Lemma 5-3).

Corollary 1. Let f be differentiable on an interval. Any zeros of f within the interval are separated by zeros of the derivative.

Proof. If $x_1 < x_2$ and $f(x_1) = f(x_2) = 0$, the conditions of Lemma 5-3 are satisfied and there exists a value u such that $x_1 < u < x_2$ and $f'(u) = 0$.

As a consequence of this result we observe further that, in a given interval, a function may have at most one more zero than its derivative. From this fact there follows a familiar result:

Corollary 2. A polynomial of degree n can have no more than n distinct real zeros.

The proof is left as an exercise (Exercises 5-3, No. 1).

Example 5-3a. (i) Let us apply Corollary 1 to the zeros of $f(x) = x^3 - 3x + 1$. We know that $f'(x) = 3x^2 - 3$ has zeros at $x = 1$ and $x = -1$. It follows that f may have as many as three zeros. We observe that $f(-1) = 3$ and $f(1) = -1$. By the Intermediate Value Theorem we conclude that there is a zero of f between -1 and 1 . Clearly we can make $f(x)$ negative for sufficiently large negative values and positive for sufficiently large positive values. It follows that f has a zero for $x < -1$ and another for $x > 1$. Specifically, we have $f(-2) = -1$ and $f(2) = 3$, so that there is one zero between -2 and -1 and another between 1 and 2 .

(ii) The function $f(x) = x^3 + 3x + 1$ has the derivative $f'(x) = 3x^2 + 3$ which is always positive. Since the derivative is always positive f can have at most one zero. Observing that $f(-1) = -3$ and $f(0) = 1$ we see that a zero exists and lies between $x = -1$ and $x = 0$.

(iii) The function $f(x) = x^4 - 4x^3 - 8x^2 + 64$ has the derivative $f'(x) = 4x^3 - 12x^2 - 16x$ which has zeros at $x = -1$, $x = 0$ and $x = 4$. Thus f may have as many as four zeros. We have $f(-1) = 61$, $f(0) = 64$ and $f(4) = -64$. It follows that f has a local maximum at $x = 0$ and from Theorem 5-2b that f is increasing on the interval $(-1, 0)$; consequently, there is no zero between -1 and 0 . Further, from

$f(-2) = 80$ we see that f has a local minimum at $x = -1$; hence that f is decreasing for $x < -1$, and there is no zero to the left of -1 . Finally, we observe that $f(6) = 208$ so that $f(x)$ has precisely two zeros, one, between 0 and 4, another between 4 and 6.

Because the curve of Figures 5-3a, b is drawn overly simply it would be easy to leap to the conclusion that the Law of the Mean is geometrically identical with Rolle's Theorem in a rotated coordinate frame for which the x -axis is parallel to the chord, but this is not so. In such a coordinate frame we may lose the property that the graph may be represented by a function. Thus, in Figure 5-3c the perpendicular to the chord from the point $(x, f(x))$ of the arc intersects the graph in more than one point. This difficulty is surmounted by expressing the distance of a point $(x, f(x))$ on the arc from the chord (the ordinate in the rotated coordinate frame) as $d(x)$, a function of the number x in the domain of f , rather than attempting to express it in terms of distance along the chord (the abscissa in the rotated coordinate frame).

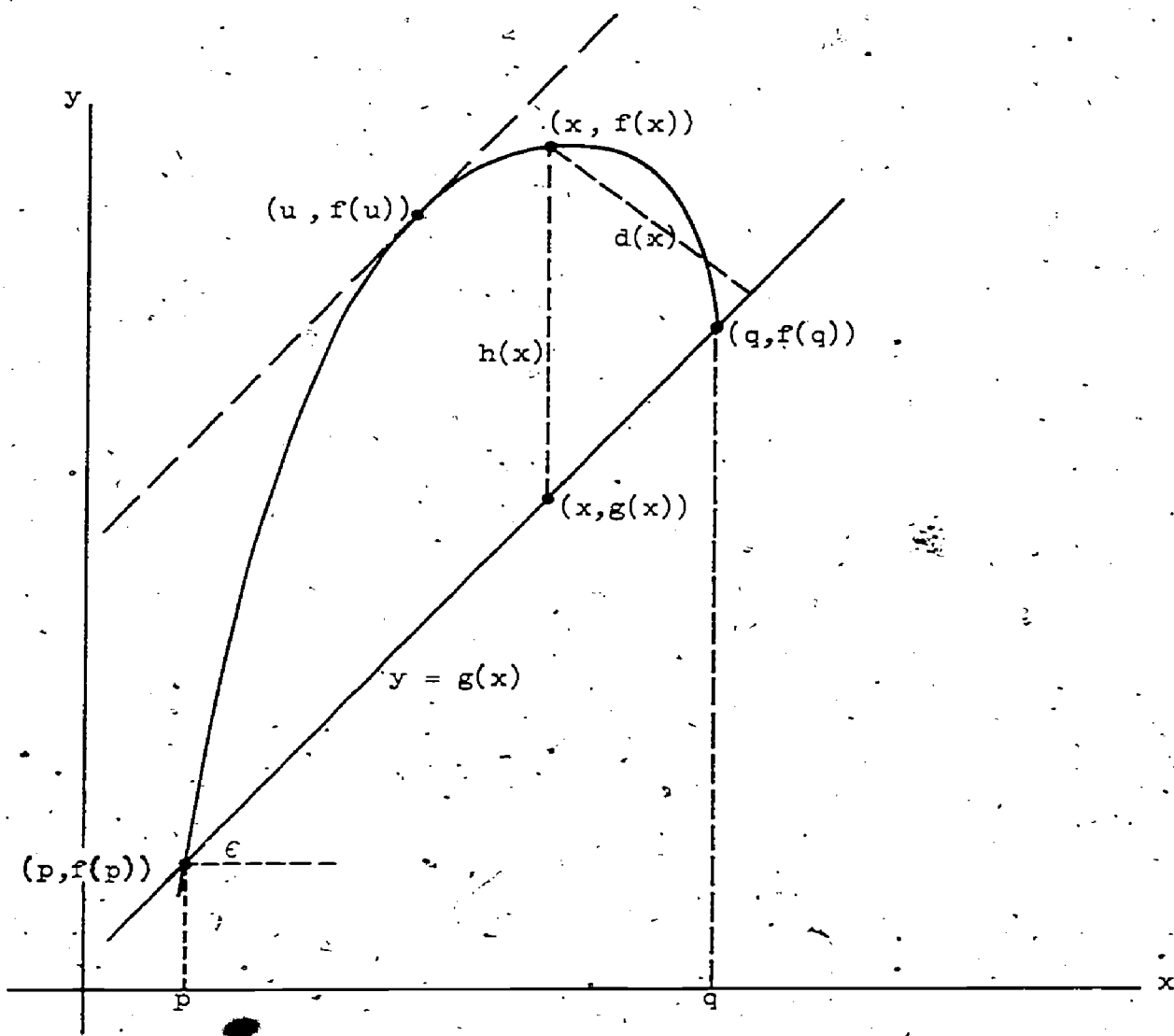


Figure 5-3c

THEOREM 5-3. (Law of the Mean). If f is continuous on the closed interval $[p, q]$ and differentiable on the open interval (p, q) then there is at least one point u in the open interval where

$$(1) \quad f'(u) = \frac{f(q) - f(p)}{q - p}$$

Proof. It is more convenient to deal not with the perpendicular distance $d(x)$ of the point $(x, f(x))$ on the arc from the line joining its two endpoints (see Figure 5-3c) but with the vertical height $h(x)$ of the point above the chord. The two quantities are proportional: $d(x) = h(x) \cos \theta$ where θ is the acute angle made by the chord with any horizontal line. An extreme value of the one is, therefore, an extreme value of the other. The equation of the straight line joining the points $(p, f(p))$ and $(q, f(q))$ is

$$(2) \quad y = g(x) = f(p) + (x - p) \frac{f(q) - f(p)}{q - p}$$

It follows for any point x in (p, q) that the height $h(x)$ of $(x, f(x))$ above the chord is given by

$$(3) \quad h(x) = f(x) - g(x) = f(x) - f(p) - (x - p) \frac{f(q) - f(p)}{q - p}$$

From this equation it follows straightforwardly that $h(x)$ satisfies the conditions of Rolle's Theorem (Lemma 5-3) on $[p, q]$. First, as you may verify directly, $h(p) = h(q) = 0$. Next observe that $h(x) = f(x) - g(x)$ is the sum of $f(x)$ and a linear function; since both terms of this sum are differentiable on the open interval (p, q) and continuous on the closed interval $[p, q]$ it follows from the theorems on the derivative of a linear combination and on the sums of continuous functions (Theorems 4-2a and 3-6a) that h also is differentiable on the open interval and continuous on the closed interval. From Rolle's Theorem, we conclude that for some value u in (p, q)

$$h'(u) = f'(u) - g'(u) = 0,$$

or, from Equation (3) for $h(x)$ above,

$$f'(u) - \frac{f(q) - f(p)}{q - p} = 0.$$

(ii) Linear interpolation.

Linear interpolation is a useful method of approximation to the values of a function in an interval when the endpoint values are known. If bounds on the range of the derivative can be obtained, the Law of the Mean gives a way of estimating the error of approximation.

Geometrically, linear interpolation consists of replacing the arc of the graph of f on (p, q) by the chord joining the endpoints. Thus, on (p, q) we approximate $f(x)$ by the linear function $g(x)$ given in Equation (2). The error of the approximation $g(x) - f(x) = -h(x)$ is given by Equation (3). For our purposes it is convenient to recast Equation (3) in the form

$$g(x) - f(x) = (x - p) \left(\frac{f(p) - f(q)}{p - q} - \frac{f(x) - f(p)}{x - p} \right).$$

Now, by the Law of the Mean

$$(4) \quad g(x) - f(x) = (x - p)[f'(u_2) - f'(u_1)]$$

where $p < u_1 < x < q$, $p < u_2 < q$. If the derivative is bounded in (p, q) , say $|f'(z)| \leq M_1$ for z in (p, q) , then from Equation (4)

$$|g(x) - f(x)| \leq |x - p| (|f'(u_2)| + |f'(u_1)|)$$

whence

$$(5) \quad |g(x) - f(x)| \leq 2M_1 |x - p|$$

Example 5-3b. Let us estimate $\sqrt{10}$ by linear interpolation for the function $f : \mathbb{R} \rightarrow \sqrt{x}$. Since $3 < \sqrt{10} < 4$ we take $p = 9$ and $q = 16$ in Equation (2) and obtain $g(10) = \frac{22}{7}$ as our estimate for $\sqrt{10}$. On the interval $(9, 16)$ we have

$$f'(x) = \frac{1}{2\sqrt{x}} < \frac{1}{2\sqrt{9}} \leq \frac{1}{6}.$$

Entering this bound in (5) we obtain

$$\left| \frac{22}{7} - \sqrt{10} \right| \leq \frac{1}{3}.$$

We observe, however, that

$$\left(\frac{22}{7}\right)^2 = \frac{484}{49} = 10 - \frac{6}{49}$$

and we suspect that our estimate of error is rather crude.

If on the interval (p, q) f' has a derivative f'' , the second derivative of f , we may apply the Law of the Mean again to the difference $f'(u_2) - f'(u_1)$ in Equation (4) to obtain

$$g(x) - f(x) = (x - p)(u_2 - u_1)f''(v)$$

where v is somewhere between u_2 and u_1 . Since u_2 and u_1 are both points of (p, q) we know that the distance between the two points is less than the length of the interval:

$$|u_2 - u_1| < q - p.$$

Suppose, in addition, that we have a bound on the second derivative, $|f''(x)| \leq M_2$ on (p, q) . Then we obtain an upper estimate for the error in terms of the second derivative:

$$(6) \quad |g(x) - f(x)| \leq (x - p)(q - p)M_2.$$

Example 5-3c. Now let us use Formula (6) to obtain an estimate for the error of approximation to $\sqrt{10}$ by the linear interpolation scheme of Example 5-3b. We have

$$\begin{aligned} |f''(x)| &= \left| -\frac{1}{4x^{3/2}} \right| < \left| \frac{1}{4 \cdot 9^{3/2}} \right| \\ &\leq \frac{1}{108} \end{aligned}$$

for x in $(9, 16)$. Consequently, from (6),

$$\left| \frac{22}{7} - \sqrt{10} \right| \leq \frac{7}{108} < .065.$$

It follows that

$$3.07 < \sqrt{10} < 3.21.$$

We have obtained sharper estimates for $\sqrt{10}$ and now we can repeat the process to obtain still sharper estimates using $p = (3.07)^2$ and $q = (3.21)^2$.

Exercises 5-3

1. Prove Corollary 2 to Lemma 5-3.
2. Sketch the graphs of the functions in Example 5-3a.
3. Is the following converse of Rolle's Theorem true? If f is continuous on the closed interval $[p, q]$ and differentiable on the open interval (p, q) , and if there is at least one point u in the open interval where $f'(u) = 0$, then there are two points m and n where $p \leq m < u < n \leq q$ such that $f(m) = f(n)$.
4. Does Rolle's Theorem justify the conclusion that $\frac{dy}{dx} = 0$ for some value of x in the interval $-1 \leq x \leq 1$ for $(y + 1)^3 = x^2$?
5. Given: $f(x) = x(x - 1)(x - 2)(x - 3)(x - 4)$. Determine how many solutions $f'(x) = 0$ has and find intervals including each of these without calculating $f'(x)$.
6. Verify that Rolle's Theorem (Lemma 5-3) holds for the given function in the given interval or give a reason why it does not.
 - (a) $f: x \rightarrow x^3 + 4x^2 - 7x - 10$, $[-1, 2]$
 - (b) $f: x \rightarrow \frac{2 - x^2}{4}$, $[-1, 1]$
7. Prove that the equation

$$f(x) = x^n + px + q = 0$$
 cannot have more than two real solutions for an even integer n nor more than three real solutions for an odd n . Use Rolle's Theorem. This problem can also be done without it (Exercises 3-7, No. 21).
8. A function g has a continuous second derivative on the closed interval $[a, b]$. The equation $g(x) = 0$ has three different solutions in the open interval (a, b) . Show that the equation $g''(x) = 0$ has at least one solution in the open interval (a, b) .
9. Show that the conclusion of the Law of the Mean does not follow for $f(x) = \tan x$ in the interval $1.5 < x < 1.6$.

10. For each of the following functions show that the Law of the Mean fails to hold on the interval $[-a, a]$ if $a > 0$. Explain why the theorem fails.
- (a) $f : x \longrightarrow |x|$
- (b) $f : x \longrightarrow \frac{1}{x}$
11. Show that the equation $x^5 + x^3 - x - 2 = 0$ has exactly one solution in the open interval $(1, 2)$.
12. Show that $x^2 = x \sin x + \cos x$ for exactly two real values of x .
13. Find a number that can be chosen as the number u in the Law of the Mean for the given function and interval.
- (a) $f : x \longrightarrow \cos x, 0 \leq x \leq \frac{\pi}{2}$
- (b) $f : x \longrightarrow x^3, -1 \leq x \leq 1$
- (c) $f : x \longrightarrow x^3 - 2x^2 + 1, -1 \leq x \leq 0$
- (d) $f : x \longrightarrow \cos x + \sin x, 0 \leq x \leq 2\pi$
14. Derive each of the following inequalities by applying the Law of the Mean.
- (a) $|\sin x - \sin y| \leq |x - y|$
- (b) $\frac{x}{1+x^2} < \arctan x < x$ if $x > 0$
15. Use the Law of the Mean to approximate $\sqrt[3]{1.008}$.
16. Use the Law of the Mean to approximate $\cos 61^\circ$.
17. Show that $a \left(1 + \frac{\epsilon}{n(a^n + \epsilon)}\right) < \sqrt[n]{a^n + \epsilon} < a \left(1 + \frac{\epsilon}{na^n}\right)$ for $\epsilon > 0, a > 1, n > 1$ (n rational).
18. Using Number 17, obtain the following approximations.
- (a) $3 + \frac{1}{10} < \sqrt[3]{30} < 3 + \frac{1}{9}$
- (b) $3 + \frac{3}{5(244)} < \sqrt[5]{244} < 3 + \frac{1}{405}$
- (c) Show that the approximation $\frac{1}{2} \left(3 + \frac{3}{5(244)} + 3 + \frac{1}{405}\right)$ to $\sqrt[5]{244}$ is correct to at least 5 decimal places.

19. (a) Show that a straight line can intersect the graph of a polynomial of n -th degree at most n times.
- (b) Obtain the corresponding result for rational functions.
- (c) Could $\sin x$ or $\cos x$ be rational functions? Justify your answer.
20. Prove the intermediate value property for derivatives; namely, if f is differentiable on the closed interval $[p, q]$ then $f'(x)$ takes on every value between $f'(p)$ and $f'(q)$ in the open interval (p, q) .

5-4. Applications of the Law of the Mean.(1) Monotone functions.

In Section 5-2 we related the zeros of the derivative to the extrema of a function f . Here we wish to consider the properties of the derivative f' on those intervals where there are no interior extrema. As we know from Theorem 5-2b, on such an interval the function f must be either increasing or decreasing. If f is increasing we do not expect to find any points where the slope of the graph is negative although, as the function $f(x) = x^3$ illustrates (Example 5-2e), there may be values of x for which $f'(x) = 0$. Since we cannot exclude the possibility that $f'(x) = 0$, there is no reason to exclude the possibility that $f'(x) = 0$ on an entire interval. On such an interval $f(x)$ is constant. To include this possibility, we have introduced the concept of weakly increasing function.*

The monotone character of a function f is directly connected with the sign of the derivative f' .

THEOREM 5-4a. If f is differentiable on (a,b) and weakly increasing then $f'(x) \geq 0$ for all x in (a,b) ; if weakly decreasing, then $f'(x) \leq 0$. Conversely, if $f'(x) \geq 0$ for all x in (a,b) then f is weakly increasing on (a,b) ; if $f'(x) \leq 0$ then f is weakly decreasing.

* See Section A2-4 for a discussion of the concepts of weakly increasing function, monotone function, etc.

Proof. We consider only the increasing case here; the proof for the decreasing case is similar.

If f is weakly increasing on (a,b) , then, for x in (a,b) and h sufficiently small that $x+h$ lies in (a,b) , we have

$$\frac{f(x+h) - f(x)}{h} \geq 0.$$

whether $h > 0$ or $h < 0$. It follows by Theorem 3-4f that

$$f'(x) = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h} \geq \lim_{h \rightarrow 0} 0 \geq 0.$$

Converse: Suppose $f'(x) \geq 0$ on (a,b) . For any two values x_1, x_2 in the interval satisfying $x_1 < x_2$ we have by the Law of the Mean

$$f(x_2) - f(x_1) = f'(u)(x_2 - x_1)$$

where $x_1 < u < x_2$. Since $f'(u) \geq 0$ and $(x_2 - x_1) > 0$ we conclude that $f(x_2) - f(x_1) \geq 0$. Since $f(x_2) \geq f(x_1)$ for any pair satisfying $x_2 > x_1$, the monotone property of f is established.

Corollary 1. If $f'(x) = 0$ for all points x in the interval (a,b) then $f(x)$ is constant on (a,b) .

The proof of this corollary follows from the observation that f must be both weakly increasing and weakly decreasing, hence constant.

From the preceding corollary we can see to what extent the derivative of a function determines the function. If two functions g and f have the same derivative then $D(g - f) = Dg - Df = 0$. It follows from Corollary 1 to Theorem 5-4a that

$$g - f = c$$

where c is a constant function. Thus

$$g = f + c.$$

In words, the derivative of a function determines the function to within an additive constant. Geometrically, the graph of a function is determined by the derivative except for a possible vertical translation. This result is summarized as follows.

Corollary 2. If two functions have the same derivative on an interval they differ by a constant. Conversely, if two functions differ by a constant they have the same derivative.

Proof. The proof of the first proposition is given above. The converse is a direct consequence of the differentiation theorems of Section 4-2.

Finally, we observe that the proof of the converse proposition in Theorem 5-4a remains valid if all weak inequalities are replaced by strong ones:

THEOREM 5-4b. If $f'(x) > 0$ for all x in (a,b) , then f is increasing on (a,b) . If $f'(x) < 0$, then f is decreasing.

The converse of Theorem 5-4b is not true (Why not?). See if you independently can find a condition on the derivative of f which is equivalent to requiring that f is strongly-monotone (Exercises 5-4, No. 14).

(ii) The reversal of sign test for an extremum.

In the examples of Section 5-2 we found the local behavior of f near a point u where $f'(u) = 0$, namely, whether $f(u)$ is a local maximum, minimum, or neither. For this purpose we compared $f(u)$ with $f(p)$ and $f(q)$ where $p < u < q$ and p, q are chosen as next adjacent zeros of f' or as endpoints of the interval. Since the values of f in any neighborhood of u are sufficient to establish a local property, we are led to seek a criterion which does not depend upon the values of f at distant points, such as other zeros of f' .

THEOREM 5-4c: Let f be differentiable on a neighborhood of a point a for which $f'(a) = 0$. If $f'(x)$ reverses sign at a then $f(a)$ is an extremum on the neighborhood. Specifically, if $f'(x) < 0$ when $x < a$ and $f'(x) > 0$ when $x > a$ then $f(a)$ is a minimum; if $f'(x) > 0$ when $x < a$ and $f'(x) < 0$ when $x > a$ then $f(a)$ is a maximum. (See Figure 5-4a.)

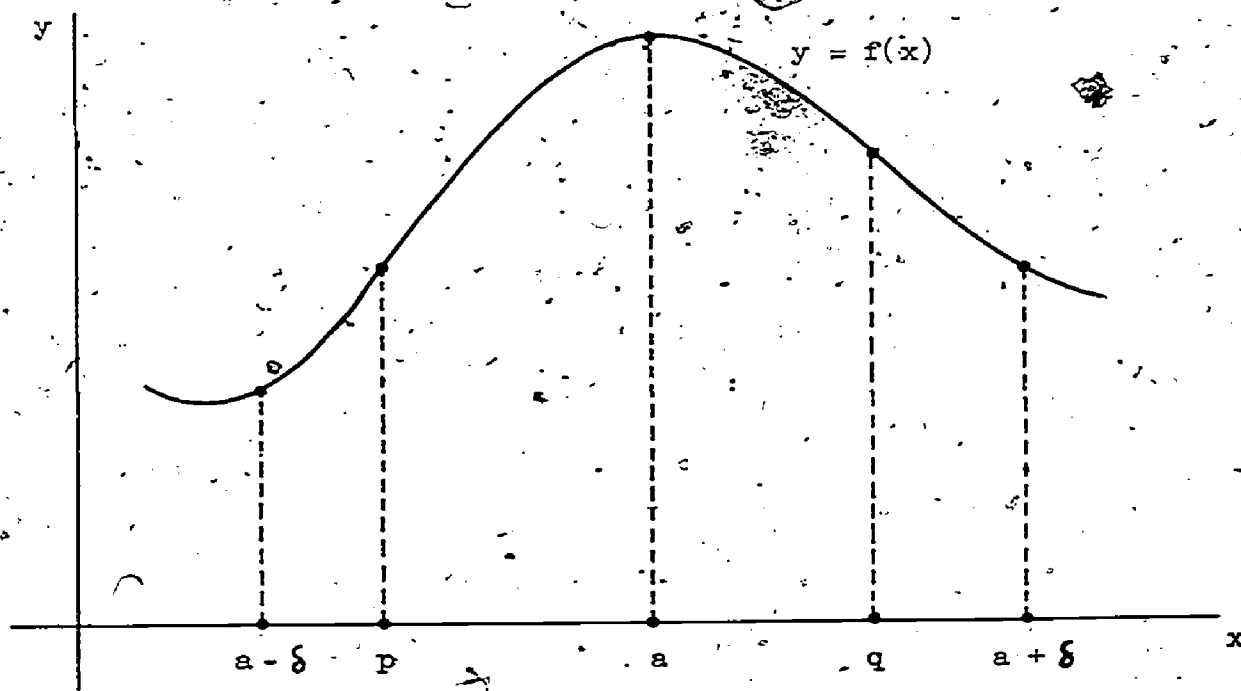


Figure 5-4a

Proof. Here we prove the criterion for a maximum. Let x be a point of the neighborhood other than a . By the Law of the Mean we have

$$f(a) - f(x) = f'(u)(a - x)$$

where u lies between a and x . Whether $x > a$ or $x < a$ it follows from the conditions of the theorem that $f'(u)(a - x) > 0$. We conclude that $f(a) > f(x)$ for all x in the neighborhood.

The conclusion of the theorem remains valid if the inequalities governing the sign of the derivative are weak, but we seldom need the theorem in that form. (Exercises 5-4, No. 11.) The proof parallels that of Theorem 5-4c. If we have strong inequalities, however, we have a sharper result.

Corollary. Under the conditions of Theorem 5-4c, the extremum of f at $x = a$ is isolated, that is, in the deleted neighborhood of a , we have $f(x) \neq f(a)$.

Proof. For the given neighborhood of Theorem 5-4c we have already proved when $f(a)$ is a maximum that there is a strong inequality $f(x) < f(a)$ for $x \neq a$. A similar argument establishes the result $f(x) > f(a)$ for $x \neq a$ when $f(a)$ is a minimum.

To complement the preceding theorem we need to know when $f(a)$ is not a local extremum.

THEOREM 5-4d. Let f be differentiable on a neighborhood of the point a for which $f'(a) = 0$. If $f'(x)$ has constant non-zero sign throughout a deleted neighborhood of a then $f(a)$ is not a local extremum.

Proof. Suppose that $f'(x)$ is positive on a deleted neighborhood of a . Let $h > 0$ be any value smaller than the radius of the neighborhood. From Theorems 5-4b and 5-2b it follows that f is increasing on both closed intervals $[a - h, a]$ and $[a, a + h]$. Consequently,

$$f(a - h) < f(a) < f(a + h).$$

It follows that a can be neither a local minimum nor a local maximum.

Example 5-4a. If $f(x) = 3x^7 - 5x^4 + \sqrt{2}x^2 + 93$ then $f'(x) = 21x^6 - 20x^3 + 2\sqrt{2}x$. Observing that $f'(0) = 0$, we ask whether $f(0)$ is a local minimum or maximum. To find the other zeros of $f'(x)$ we would have to solve the fifth degree equation $21x^5 - 20x^2 + 2\sqrt{2} = 0$. It seems preferable to test for reversal of sign. Writing

$$f'(x) = x(2\sqrt{2} - 20x^2 + 21x^5)$$

we observe that for x near zero the factor in parentheses is close to $2\sqrt{2}$, hence positive. It follows that for sufficiently small x , $f'(x)$ changes sign with x ; if $x < 0$ then $f'(x) < 0$ and if $x > 0$ then $f'(x) > 0$. We conclude that $f(0) = 93$ is a local minimum.

Example 5-4b. If $f(x) = x^8 - 7x^5 + \pi x^3 - 1024$ then

$f'(x) = 8x^7 - 35x^4 + 3\pi x^2$ and $f'(0) = 0$. We test for reversal of sign. Writing

$$f'(x) = x^2(3\pi - 35x^2 + 8x^5)$$

we see that the factor in parentheses is positive when x is sufficiently small. It follows that $f'(x)$ is positive whether $x > 0$ or $x < 0$. We conclude by Theorem 5-4d that $f(0)$ is neither a local maximum nor a local minimum, but that f is an increasing function in the neighborhood of $x = 0$.

The reversal-of-sign criterion is especially valuable when there is only one zero of the derivative in an interval. The criterion is then sufficient to confirm an interior global extremum. (See Exercises 5-4. No. 12.)

Example 5-4c. A cord of length L has a small ring attached to one end; the other end is first passed through the ring so as to form a loop, then fastened to a weight. If the loop is passed over two horizontal pegs a distance M apart, jutting out of a wall at the same level, at what height h below the level of the pegs will the weight come to rest (Figure 5-4b)? (It is assumed that L is considerably longer than $2M$, that the height of the pegs above the ground is greater than L , and that friction plays no role.)

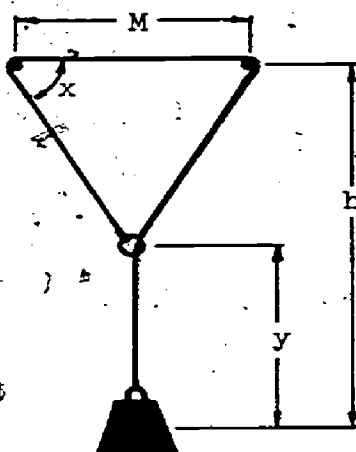


Figure 5-4b

We denote the distance between the level of the weight and the level of the pegs by h , and the distance between the levels of the weight and ring by y . If x denotes the angle formed by the cord at either peg, then we have

$$h = y + \left(\frac{M}{2}\right) \tan x, \quad L = y + M + 2\left(\frac{M}{2}\right) \sec x,$$

where we may assume $0 \leq x < \frac{\pi}{2}$. We eliminate y between these two relations, and obtain a function f given by the equation

$$f(x) = h = \left(\frac{M}{2}\right) \tan x + L - M - M \sec x,$$

whose maximum we seek. The derivative

$$\begin{aligned} f'(x) &= \left(\frac{M}{2}\right) \sec^2 x - M \sec x \tan x \\ &= \left(\frac{M}{2}\right) \sec x (\sec x - 2 \tan x) \end{aligned}$$

is zero only for $\sec x = 2 \tan x$, or $\sin x = \frac{1}{2}$ --that is, for $x = \frac{\pi}{6}$.
From the rearranged form

$$f'(x) = \left(\frac{M}{2}\right) \sec^2 x (1 - 2 \sin x).$$

we see that $f'(x) > 0$ for $0 \leq x < \frac{\pi}{6}$ and $f'(x) < 0$ for $\frac{\pi}{6} < x < \frac{\pi}{3}$, so that

$$\begin{aligned} f\left(\frac{\pi}{6}\right) &= \left(\frac{M}{2}\right) \tan \frac{\pi}{6} + L - M - M \sec \frac{\pi}{6} \\ &= L - \left(\frac{M}{2}\right)(2 + \sqrt{3}) \end{aligned}$$

is a local maximum value of f . It would be tedious, but not difficult, to show that no greater value of f is attained at the endpoints of the interval of physically possible values of the angle x --when the ring is at the level of the pegs or when the ring is at the level of the weight. We observe, however, that there is no other zero of f' in the extended domain $0 \leq x < \frac{\pi}{2}$ so that without further test we know that $f\left(\frac{\pi}{6}\right)$ is the overall maximum and describes the equilibrium position.

Exercises 5-4

1. On what intervals is the function

$$f : x \longrightarrow \frac{x^2 - 3}{x - 2}$$

strongly monotone? Use Theorem 5-4c to characterize all extrema.

2. Locate all intervals on which the function

$$f(x) = 44 + 4x - 13x^2 + 18x^3 - 9x^4$$

is increasing; decreasing. (Compare with your answer to Exercises 5-2a, No. 3.)

3. For each of the following functions find all points a for which $f'(a) = 0$. Examine the sign of f' and determine those intervals on which f is strongly monotone.

(a) $f(x) = \frac{x}{1+x^4}$

(b) $f(x) = (1-x)^4$

(c) $f(x) = (1-x)^5$

(d) $f(x) = \frac{x^2+10}{x^2-5}$

4. If p and q are integers and

$$f(x) = (x-1)^p(x+1)^q, \quad (p \geq 2, q \geq 2)$$

find the extrema of f , for the following cases:

(a) p and q are both even.

(b) p is even and q is odd.

(c) p is odd and q is even.

(d) p and q are both odd.

5. If p , q , and r are positive integers, and $a < b < c$, discuss the graph of the function

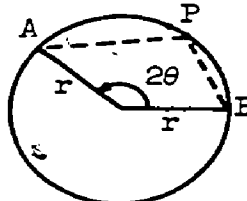
$$f: x \longrightarrow (x-a)^p(x-b)^q(x-c)^r.$$

Discuss some special cases as in Number 4.

6. A tank is to have a given volume V and is to be made in the form of a right circular cylinder with hemispherical ends. The material for the ends costs twice as much per square foot as that for the cylindrical part. Find the most economical dimensions.

7. Find the length of the longest rod which can be carried horizontally around a corner from a corridor 10 ft. wide into one 5 ft. wide.

8. Find a point P on the arc AB such that the sum of the lengths of the chords AP and BP is a maximum ($\theta \leq \frac{\pi}{2}$).



9. Show how to construct a line, if possible, which passes through the point $(5,8)$ such that the area of the triangle formed in the first quadrant is a positive number a . For what values of a is it impossible to construct such a triangle?
10. Find a point on the altitude of an isosceles triangle such that the sum of its distances from the vertices is the smallest possible.
11. Let f be differentiable on a neighborhood of a point a for which $f'(a) = 0$. If $f'(x) \leq 0$ when $x < a$ and $f'(x) \geq 0$ when $x > a$, then $f(a)$ is a minimum. If $f'(x) \geq 0$ when $x < a$ and $f'(x) \leq 0$ when $x > a$ then $f(a)$ is a maximum. Give a proof.
12. Let f be continuous on the closed interval $[a,b]$ and differentiable on the open interval (a,b) . Suppose u is the one point in (a,b) where $f'(u) = 0$. Prove that if $f'(x)$ reverses sign in a neighborhood of u then $f(u)$ is the global extremum of f on $[a,b]$ appropriate to the sense of reversal.
13. Given a function f such that $f(1) = f(2) = 4$, and such that $f''(x)$ exists and is positive throughout the interval $1 \leq x \leq 3$.
- What can you conclude about $f'(2.5)$?
 - Prove your statement, stating whatever theorems you use in your proof.
14. Let f be a differentiable function on (a,b) . Prove that the requirement that f be increasing is equivalent to the condition that $f'(x) \geq 0$ everywhere but that every interval contains points where $f'(x) > 0$.
15. Given that f is everywhere differentiable. If for all x such that $f'(x) \geq 0$, $f(x) \leq f(0)$, prove that $f(x) \leq f(0)$ for all $x \geq 0$.
16. A function g is such that g'' is continuous and positive in the interval (p,q) . What is the maximum number of roots of each of the equations $g(x) = 0$ and $g'(x) = 0$ in (p,q) ?
Prove your result and give some illustrative examples.
17. (a) If $f'(a) > 0$ show for values of x in a neighborhood of a that if $x > a$ then $f(x) > f(a)$, and if $x < a$ then $f(x) < f(a)$.
- (b) Give an example of a function f for which $f'(a) > 0$ but which is not increasing in any neighborhood of a , no matter how small.

5-5. Applications of the Second Derivative.

(1) Second derivative test for an extremum.

Reversal of sign of the first derivative is a sufficient condition for an extremum (Theorem 5-4c). As we survey the graph of f from left to right in a neighborhood of a , if first $f'(x)$ is negative for $x < a$, next $f'(a) = 0$, finally $f'(x)$ is positive for $x > a$, then we know that $f(a)$ is a local minimum. To guarantee a minimum at a , then, it would be sufficient to demonstrate that f' is increasing on a neighborhood of a . To establish the increasing character of f' and hence that $f(a)$ is a minimum it is sufficient by Theorem 5-4b to show that the derivative of f' , the second derivative f'' of f , is positive on a neighborhood of a . Actually it is enough to know only that $f''(a) > 0$ as we now prove.

THEOREM 5-5a. Suppose $f'(a) = 0$. If $f''(a) > 0$ then $f(a)$ is a local minimum of f . If $f''(a) < 0$ then $f(a)$ is a local maximum of f .

Proof. We consider the case $f''(a) > 0$. From the definition of limit and

$$f''(a) = \lim_{x \rightarrow a} \frac{f'(x) - f'(a)}{x - a}$$

we know that for a sufficiently small deleted neighborhood of a

$$\frac{f'(x) - f'(a)}{x - a} > 0.$$

This inequality implies for $x > a$ that $f'(x) - f'(a) > 0$, and for $x < a$ that $f'(x) - f'(a) < 0$. Observe that this proves the result of Exercises 5-4, No. 17 for f' . But this condition assures a minimum, by Theorem 5-4c.

Example 5-5a. In Example 5-4c of the weight suspended from a string looped over pegs we found for the derivative of the height function

$$f'(x) = \left(\frac{M}{2}\right) \sec^2 x (1 - 2 \sin x)$$

and thus showed that $x = \frac{\pi}{6}$ is the only zero of h' in the relevant domain. We test for an extremum and obtain,

$$f''\left(\frac{\pi}{6}\right) = -M \sec \frac{\pi}{6} = -\frac{2M}{\sqrt{3}}.$$

Since $f''\left(\frac{\pi}{6}\right) < 0$ we conclude that $f\left(\frac{\pi}{6}\right)$ is a local maximum, and since $\frac{\pi}{6}$ is the only zero of f' we conclude that the maximum is global.

Example 5-5b. Let us find the extrema of

$$f(x) = 4x^5 + 5x^4 - 20x^3 - 50x^2 - 40x.$$

We obtain the first and second derivatives:

$$\begin{aligned} f'(x) &= 20x^4 + 20x^3 - 60x^2 - 100x - 40 \\ &= 20(x+1)^3(x-2) \end{aligned}$$

and

$$\begin{aligned} f''(x) &= 20[(x+1)^3 + 3(x+1)^2(x-2)] \\ &= 20(x+1)^2(4x-5). \end{aligned}$$

The zeros of f' occur at $x = -1$ and $x = 2$. We attempt to apply the second derivative test and obtain $f''(-1) = 0$, $f''(2) = 540$. It follows that $f(2)$ is a local minimum. The criterion of Theorem 5-5a gives us no information about $f(-1)$ (but we observe that there is a reversal of sign of f' from positive to negative so that $f(-1)$ is a local maximum).

(ii) Convexity.

The sign of the second derivative corresponds to a useful geometrical property of the function. If $f''(x) > 0$ on an interval, then as x increases the graph flexes upward (Figure 5-5a). If $f''(x) < 0$, then as x increases

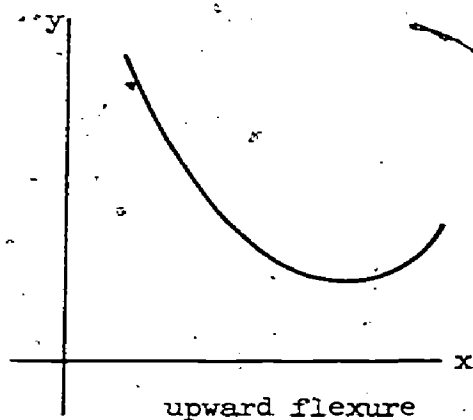


Figure 5-5a

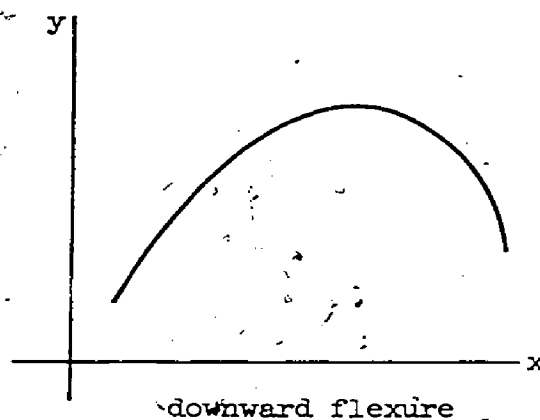


Figure 5-5b

the graph flexes downward. The geometrical concept of flexure does not require that f have first and second derivatives and we define the concept slightly more generally.

DEFINITION 5-5. The graph of the function f is said to be flexed upward on an interval of the domain if no point on any chord to the graph lies below the corresponding arc. The graph of f is said to be flexed downward if the graph of $-f$ is flexed upward. In either case, we say the function f is convex on the interval.

Note that this definition includes the limiting case of a straight line which is considered to be flexed both upward and downward.

The idea of convex function is closely linked to the geometrical concept of convex set. A set of points is convex if, for each pair of points in the set, the set contains the entire line segment joining them. Thus the interior $x^2 + y^2 < 1$ of a circle is convex; the exterior $x^2 + y^2 > 1$ is not. This idea can be expressed analytically as follows. If a and b are two points of the number line then x is a point of the closed segment joining them if and only if

$$x = \theta a + (1 - \theta)b$$

for some number θ satisfying $0 \leq \theta \leq 1$. Similarly if (a,b) and (α,β) are two points of the plane, the point (x,y) is on the closed segment joining them if and only if

$$\begin{cases} x = \theta a + (1 - \theta)\alpha \\ y = \theta b + (1 - \theta)\beta \end{cases}$$

for some θ satisfying $0 \leq \theta \leq 1$. (The verification of these assertions is left as Exercises 5-5, No. 14a.)

Figures 5-5a, b make the connection between convex sets and convex functions evident: if the graph of f is flexed upward then the set of points D above the graph is a convex set. Here the set D of points above the graph for an interval I in the domain of the function is defined by

$$D = \{(x,y) : x \in I \text{ and } y \geq f(x)\}.$$

Definition 5-5 can be expressed in analytic terms as follows. Let f be flexed upward on the interval I . Let a and b be points of I . We represent a point (x,y) of the chord joining $(a, f(a))$ to $(b, f(b))$ by

$$(\theta a + (1 - \theta)b, \theta f(a) + (1 - \theta)f(b))$$

where $0 \leq \theta \leq 1$. The statement that the graph $y = f(x)$ is flexed upward on I is equivalent to

$$f(\theta a + (1 - \theta)b) \leq \theta f(a) + (1 - \theta)f(b)$$

for all numbers a and b in I and all θ satisfying $0 \leq \theta \leq 1$.

Example 5-5c. The graph of $y = |x|$ is easily shown to be flexed upward:

$$\begin{aligned} |\theta a + (1 - \theta)b| &\leq |\theta a| + |(1 - \theta)b| \\ &\leq \theta |a| + (1 - \theta)|b| \end{aligned}$$

where $0 \leq \theta \leq 1$.

THEOREM 5-5b. Let f have a second derivative on the open interval (a, b) . The graph of f is flexed upward on (a, b) if and only if $f''(x) \geq 0$ on the interval.

The theory of convex functions and convex sets is an elegant and useful subject but to prove any of its significant results would lead us too far from our main theme. We leave to the exercises the proof of Theorem 5-5b and the derivation of some of the preliminary results of the theory. The exercises require a subtle interweaving of geometry and analysis. After exploring the problems of this section, you may wish to extend your knowledge by collateral readings.*

*For the theory of convex sets, see Yaglom, I.M. and Boltyanskii, V.G. Convex Figures. Holt, Rinehart, and Winston, 1961. For applications see Glicksman, A.M. An Introduction to Linear Programming and the Theory of Games. Wiley, New York, 1963. For the theory of convex functions at a more advanced level, see Hardy, Littlewood and Polya. Inequalities. University Press, Cambridge, England, 1953.

Exercises 5-5

1. For each of the following functions, locate and characterize all extrema and state the intervals on which the function is increasing (decreasing). On what intervals is the graph flexed upward? downward?

(a) $f : x \longrightarrow x^2 + x^{-2}$

(b) $f : x \longrightarrow \frac{2}{x^2} + \frac{1}{x}$

(c) $f : x \longrightarrow x^{1/2} + x^{-1/2}$

(d) $f : x \longrightarrow x^3 + ax^2 + bx + c$

2. Show that the graph of the function

$$f : x \longrightarrow 3 \sin 2x + 5 \cos 2x$$

is flexed upward when $f(x) < 0$ and flexed downward when $f(x) > 0$.

3. Find and characterize the extrema of the function

$$f : x \longrightarrow x \sin x + \cos x$$

on the closed interval $[0, \pi]$. On what intervals is the graph of the function flexed downward? upward?

4. (a) Assume that the function $f(x)$ has a local maximum at $x = a$, where $f'(a) = 0$, and $f''(a) \neq 0$. Determine conditions on a function g , assumed differentiable, such that $gf(x)$ also has a local maximum at $x = a$.

(b) What are the corresponding conditions for a minimum?

5. Use Theorem 5-5a to locate and classify all extrema of the function

$$f : x \longrightarrow \sin x(1 + \cos x)$$

on the closed interval $[-\pi, \pi]$. On what intervals is the graph of the function flexed downward? upward?

6. Let $f(x) = x - \sin x$. Does f have any extrema? Justify your answer.

7. Let $f(x) = \frac{x}{\sin x}$. Does f have any extrema in the open interval $(0, \frac{\pi}{2})$? Justify your answer.

8. At what point of the positive x -axis is the angle subtended by the two points $(0, 3)$ and $(4, 7)$ greatest?

9. Suppose that $f^{(1)}(a) = f^{(2)}(a) = \dots = f^{(n-1)}(a) = 0$ but that $f^{(n)}(a) \neq 0$. Determine whether $f(a)$ is a local extremum, and if it is, which kind. (Hint: consider separately the cases n even and n odd.)
10. Prove that if the graph of f is flexed downward on an interval I , then the set D of points under the graph is a convex set.
11. Prove that a necessary and sufficient condition that the graph of f be flexed downward on an interval I is that for each point a in I , the slope of the chord joining a point $(x, f(x))$ to the fixed point $(a, f(a))$ is a weakly decreasing function of x on I . (See Figure 5-5a.)
12. (a) Let f be differentiable and its graph be flexed downward on an interval I . Prove that the function

$$\phi(x) = \begin{cases} \frac{f(x) - f(a)}{x - a}, & x \neq a \\ f'(a), & x = a \end{cases}$$

is weakly decreasing, where the fixed point a is any interior point of I .

- (b) From the result of (a), prove that a necessary and sufficient condition that the graph of f be flexed downward on I is that f' be weakly decreasing.

13. Prove Theorem 5-5b.

14. (a) Let x and y be two points on an interval I in the domain of a function f . Show that a point is on the chord joining the points $(x, f(x))$ and $(y, f(y))$ on the graph of f if, and only if, its coordinates are

$$(\theta x + (1 - \theta)y, \theta f(x) + (1 - \theta)f(y))$$

for some θ such that $0 \leq \theta \leq 1$.

- (b) Show that Definition 5-5 asserts that f is flexed upward on I if, and only if, for all x and y in I and all θ such that $0 \leq \theta \leq 1$,

$$f(\theta x + (1 - \theta)y) \leq \theta f(x) + (1 - \theta)f(y).$$

- (c) Use (b) to show that the graphs of the following functions are flexed upward.

$$(i) \quad f : x \longrightarrow ax + b$$

$$(ii) \quad f : x \longrightarrow x^2$$

$$(iii) \quad f : x \longrightarrow -\sqrt{x}$$

15. (a) Derive the following property of convex functions. If the graph of f is flexed downward on an interval I , then for all points a , b in I and any positive numbers p , q

$$f\left(\frac{pa + qb}{p + q}\right) \geq \frac{pf(a) + qf(b)}{p + q}$$

In words, the function value of a weighted average is not less than the weighted average of the function values.

- (b) Prove that this property is sufficient for downward flexure.

16. Prove that if f is continuous, then a necessary and sufficient condition for its graph to be flexed downward is that

$$f\left(\frac{a + b}{2}\right) \geq \frac{f(a) + f(b)}{2}$$

17. The graph of a function f is flexed downward and is positive for all x . Show that f is a constant function. Do not assume f' exists.
18. Under what circumstances will the graph of a function f and its inverse both be flexed downward? one flexed downward and the other upward? Answer this question both with and without calculus.
19. If either of $D^2x F(x)$ or $D^2F\left(\frac{1}{x}\right)$ is of one sign for $x > 0$, show that the other one has the same sign. Interpret geometrically and illustrate by several examples.
20. If $F(x)$ is flexed upward and $F(a) = F(b) = F(c)$ where $a < b < c$, show that $F(x)$ is constant in (a, c) .
21. Show that an increasing convex function of a convex function is convex.

22. (a) Let a , b , and c be points in I such that $a < b < c$, and suppose that the graph of f is flexed upward in I . Show that

$$f(b) \leq \frac{c-b}{c-a} f(a) + \frac{b-a}{c-a} f(c).$$

(Hint: use the result of Number 15.);

hence,

$$f(a) \geq \frac{c-b}{c-a} f(b) - \frac{b-a}{c-a} f(c),$$

$$f(c) \geq \frac{c-a}{b-a} f(b) - \frac{c-b}{b-a} f(a).$$

- Λ(b) If the graph of F is flexed upward in a closed interval, show that F is bounded in the interval.
- Λ(c) Show by a counter-example that the result in (b) is not valid for an open interval.
- Λ23. (a) If the graph of F is flexed upward in an open interval, show that F is continuous in the interval.
- (b) Show by a counter-example that the result in (a) is not valid for a closed interval.
- Λ24. If the graph of f is flexed upward in an interval, then f possesses left- and right-sided derivatives at each interior point of the interval (Exercises 4-2a, No. 7).

5-6. Constrained Extreme Value Problems.

In this section we treat the problem of finding the extreme value of a function of more than one variable subject to certain side conditions or constraints.

Example 5-6a. Problem: to find the point or points nearest to the origin on the parabola

$$4y^2 - 16y - 12x - 75 = 0.$$

We wish to minimize the distance of (x,y) from $(0,0)$, namely

$$d(x,y) = \sqrt{x^2 + y^2}$$

subject to the constraint,

$$g(x,y) = 4y^2 - 16y - 12x - 75 = 0.$$

However, to avoid dealing with square roots, we may equivalently minimize

$$f(x,y) = x^2 + y^2 = [d(x,y)]^2.$$

Let us assume that x is implicitly defined by the constraint as a differentiable function of y , that is, $x = \phi(y)$ where

$$g(\phi(y), y) = 0.$$

Now set $u = f(\phi(y), y)$, $v = g(\phi(y), y)$. Differentiation with respect to y yields

$$\frac{dv}{dy} = 8y - 16 - 12 \frac{dx}{dy} = 0$$

because $g(x,y)$ is constrained to be a constant, and also

$$\frac{du}{dy} = 2x \frac{dx}{dy} + 2y.$$

Elimination of $\frac{dx}{dy}$ gives

$$\frac{du}{dy} = \frac{2}{3} (2xy - 4x + 3y).$$

For an extremum, we must have $\frac{du}{dy} = 0$. Simultaneous solution of the two equations $\frac{du}{dy} = 0$ and $g(x,y) = 0$, or

$$x = \frac{3y}{4 - 2y}, \quad 4y^2 - 16y - 12x = 75$$

shows that the y-coordinate of each extreme point must satisfy

$$8y^3 - 48y^2 - 50y + 300 = 0.$$

We factor: $0 = 8y^2(y - 6) - 50(y - 6) = 8(y^2 - \frac{25}{4})(y - 6)$. The roots $y = -\frac{5}{2}, \frac{5}{2}, 6$ have corresponding x-coordinates $x = -\frac{5}{6}, -\frac{15}{2}, -\frac{9}{4}$.

Computation of the distances from the origin for the three possible extrema shows that $(-\frac{5}{6}, -\frac{5}{2})$ is the nearest point, at a distance of $\frac{5}{6}\sqrt{10}$ (Figure 5-6a).

The general idea of a constrained extreme value problem of this type is easily visualized geometrically. We wish to find the extrema of $f(x,y)$ subject to the constraint $g(x,y) = 0$. We may think of $f(x,y)$ in a three-dimensional frame of reference as the height above the x,y-plane of a point on the surface $z = f(x,y)$. The equation $g(x,y) = 0$ may be thought of as the equation of a cylinder (whose elements are parallel to the z-axis) which meets the surface along a curve. The extreme value problem is to determine the high and low points of this curve. The general situation is depicted in Figure 5-6b. In Figure 5-6c we show the picture for Example 5-6a.

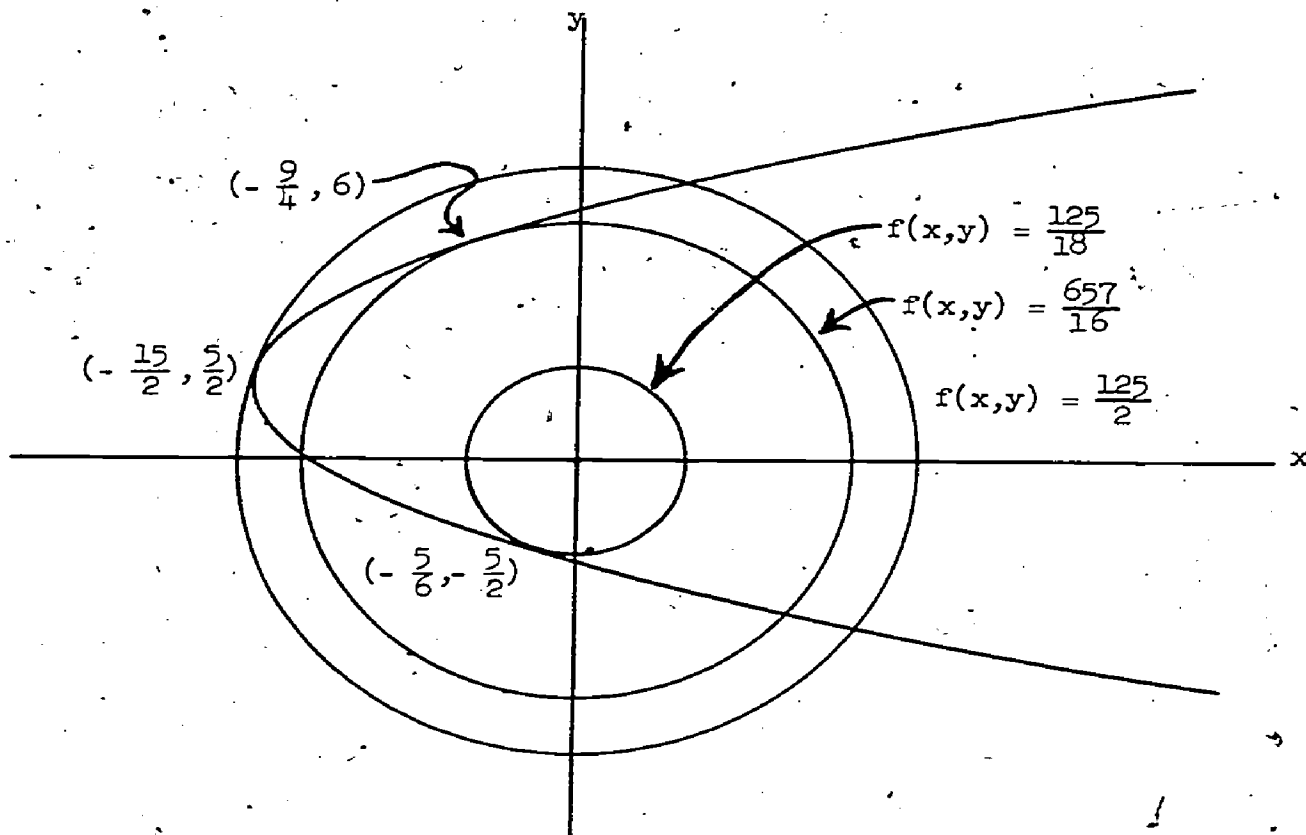


Figure 5-6a

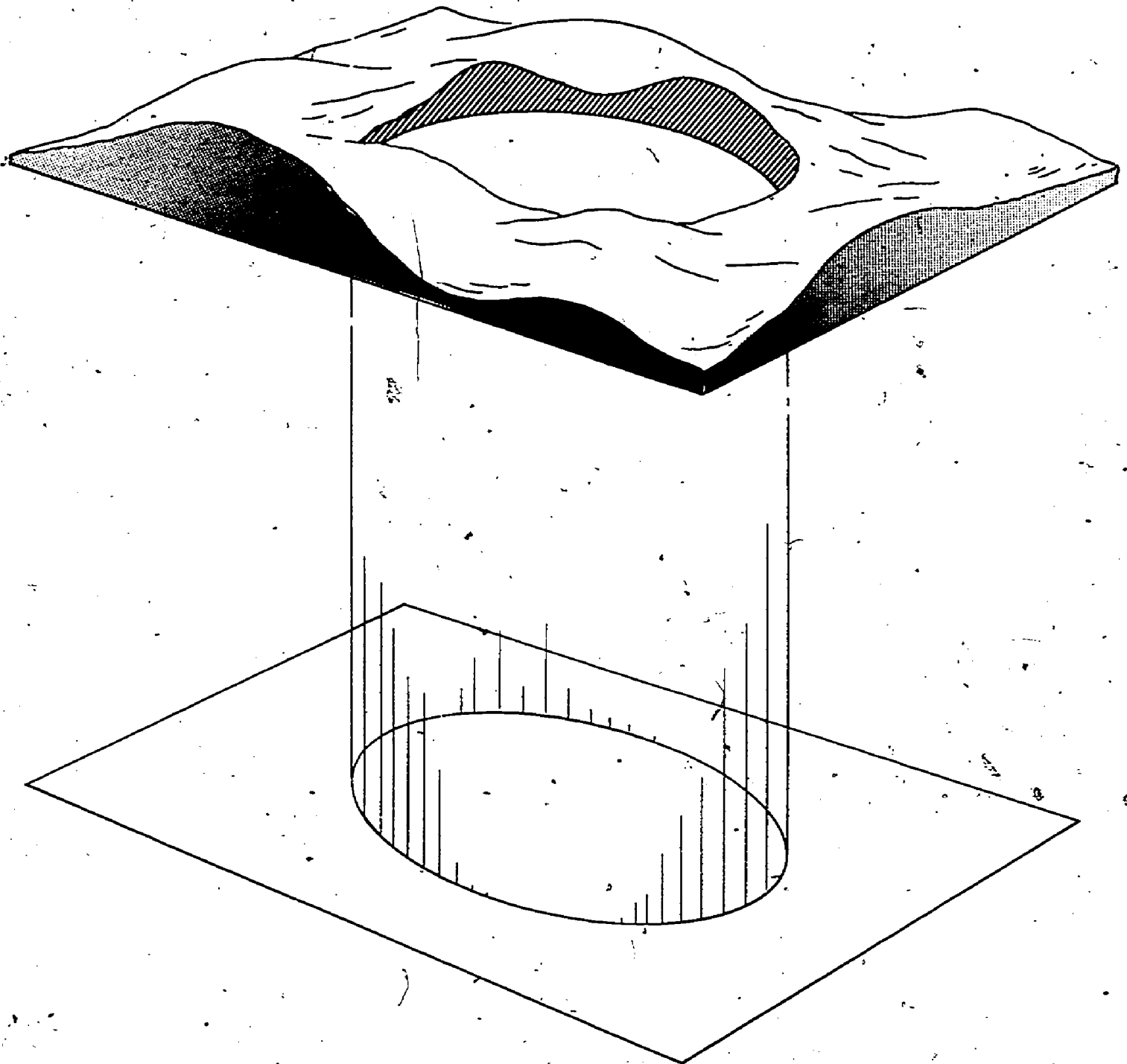


Figure 5-6b

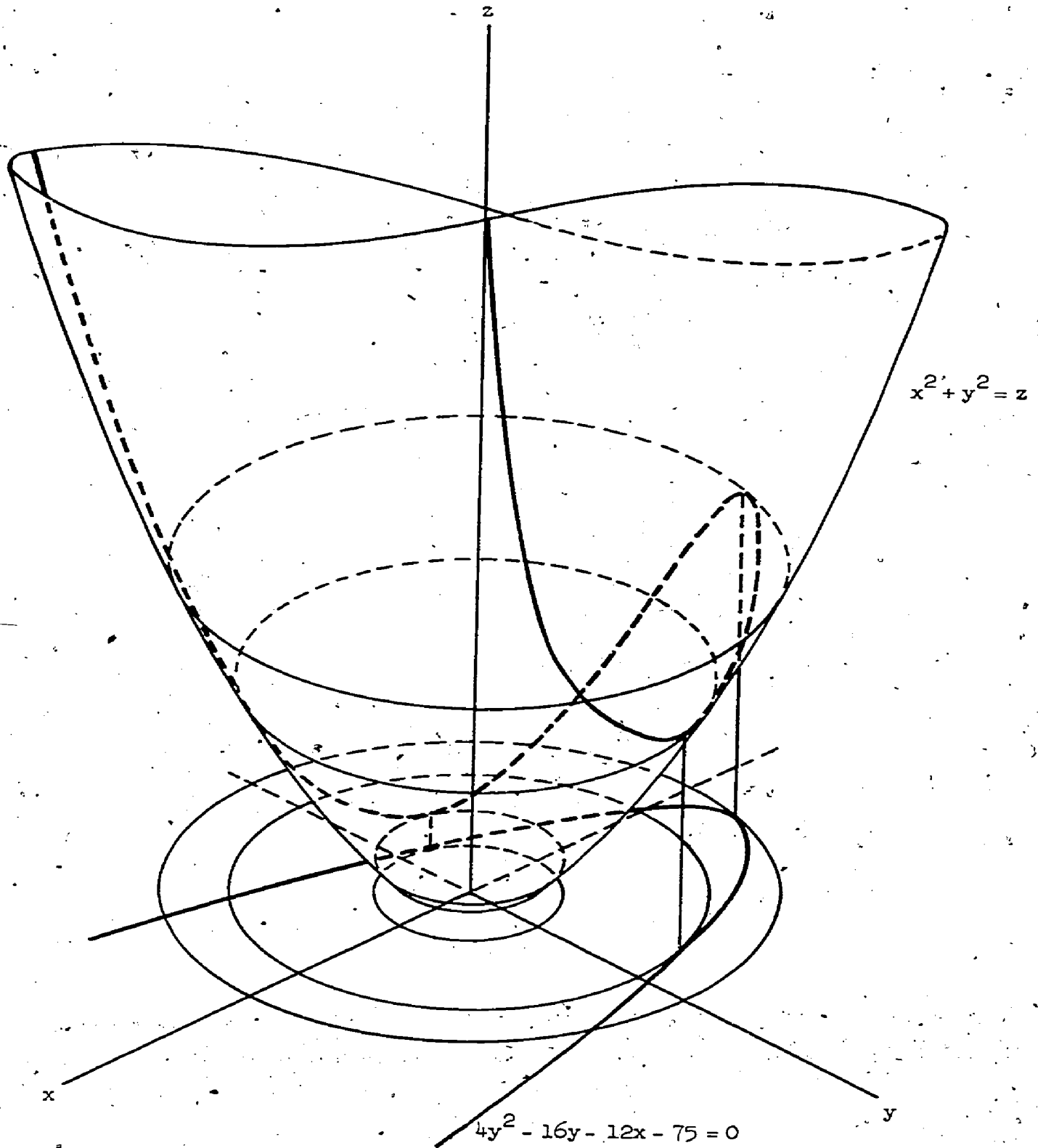


Figure 5-6c

A natural approach to the constrained extreme value problem, i.e., find the extrema of $f(x,y)$ subject to $g(x,y) = 0$, is to solve the equation $g(x,y) = 0$ to express one variable in terms of the other and to substitute this expression in $f(x,y)$. The problem is then reduced to the ordinary one of extremizing a function of a single variable. As we have indicated in Section 4-8, however, such explicit representations are often difficult or impossible to obtain (furthermore, if obtained they may not be particularly useful).

We consider the problem of maximizing a function of three variables subject to two constraints in order to show how the techniques can be extended to more complicated cases.

Example 5-6b. In Example 5-4c we might have introduced the quantity z which represents the difference in level between the ring and the pegs (see Figure 5-4b, p. 201.) and reformulated the problem as the problem of finding the maximum of

$$h = y + z$$

if y and z are subject to the constraints

$$M = 2z \cot x, \quad L = M + y + 2z \csc x.$$

Of course, we could eliminate z and y and obtain the same problem as before but we need not do so. This time, we treat x and y as functions of z and obtain

$$h' = y' + 1$$

and from the constraining conditions we obtain

$$2 \cot x - 2zx' \csc^2 x = 0$$

$$y' + 2 \csc x - 2x'z \csc x \cot x = 0,$$

since L, M are constants. From the first of the equations we obtain

$$zx' = \cos^2 x \sin x$$

and substituting for zx' in the equation for y' we obtain

$$y' + 2 \csc x (1 - \cos^2 x) = y' + 2 \sin x = 0.$$

The condition for a maximum then becomes

$$h' = 1 - 2 \sin x = 0$$

which yields the same condition as before.

There is no special merit for preferring the present treatment of this problem to the earlier one. Our only purpose is to show that elimination is not necessary should it be inconvenient or difficult.

The next example is cautionary.

Example 5-6c. Given 15 yards of fencing, I decide to plant one square and one circular flower bed, and to surround them with the fencing. What should be the dimensions of the two fences so as to contain flower beds of greatest possible area?

We express the problem in terms of the side s of the square bed and the radius r of the circle, leaving them in the implicit functional relationship

$$4s + 2\pi r = 15.$$

Subject to this constraint, s and r are to be chosen so as to maximize

$$A = s^2 + \pi r^2.$$

Denoting derivatives with respect to s by a prime, we obtain

$$4 + 2\pi r' = 0 \text{ and } A' = 2s + 2\pi r r'.$$

Setting $A' = 0$ and eliminating r' we obtain the condition for an extremum. Because of the constraint, this occurs when $8r + 2\pi r = 15$, or

$$r = \frac{15}{(8 + 2\pi)}, \quad s = \frac{30}{(8 + 2\pi)}.$$

The sum of the corresponding areas is 11.85 sq. yds. (to two places after the decimal point).

We take the usual precaution of checking the end-points of the intervals of physically possible values of r and s . If $r = 0$, then $s = \frac{15}{4}$, and we have a single square bed of area 14.06 sq. yds.; if $s = 0$, then $r = \frac{15}{2\pi}$, and we have a single round bed of area 17.90 sq. yds. Both exceed the sum

of areas found by the method of implicit differentiation, which in fact is a local minimum. It appears that the problem has no solution in the terms posed; a square and a round flower bed together will never encompass as great an area as a single round bed whose perimeter equals the total length available.

This example reveals one of the weaknesses of the method: if there should be an endpoint extremum it may be concealed by the formulation of the problem.

In the next example, we use the method to find an extremum when too few constraints are imposed by the problem. In that case we introduce an artificial constraint.

Example 5-6d. We return to the postal problem of Section 1-1, to find the linear dimensions z , x , y , of a carton which maximize the volume

$$(1) \quad V = zxy$$

subject to the constraint

$$(2) \quad 2z + 2x + y = 72, \quad (y \geq x \text{ and } y \geq z),$$

(i.e., the length y of the longest side plus the "girth" or perimeter of the cross-section perpendicular to the longest side is fixed). In the solution to this problem in Section 1-1, we indicated that to achieve a maximum we must have square ends: $z = x$. This condition served as the extra constraint necessary to complete the solution. Now we proceed in a different way.

We take as an unproven assumption that there exist dimensions z_0 , x_0 , y_0 which maximize V . If we knew y_0 we could then take $y = y_0$ and solve the problem as an ordinary constrained extreme value problem in x and z . We do not know y_0 but we can impose the artificial constraint

$$(3) \quad y = k, \quad k \text{ constant},$$

and solve the extreme value problem subject to this extra constraint. For each choice of k we obtain the largest possible value V ; this association will ordinarily be expressible as a function of k .

$$f: k \longrightarrow V.$$

Finally, we determine k so that $V = f(k)$ is a maximum.

In the above problem, then, we begin with the constrained extreme value problem, to maximize V as given by (1), subject to the constraints (2) and (3). We assume that z is a function of x . Differentiating with respect to x we obtain the condition for an extremum

$$(4) \quad \frac{dV}{dx} = y(x \frac{dz}{dx} + z) = 0$$

subject to

$$(5) \quad 2 \frac{dz}{dx} + 2 = 0$$

where we have employed condition (3) to set $\frac{dy}{dx} = 0$ in these differentiations. From (4) and (5) and the observation that $y = 0$ certainly does not maximize V , we obtain

$$z = x$$

which is the condition that the box has square ends. From this point on, the solution is conventional.

Finally, it should be remarked that the second derivative can also be computed by implicit differentiation, but that this usually entails further complication. Moreover, the most we can do with the second derivative is to infer the nature of a local extremum. Even if there is just one extremum found by this method, we cannot conclude, as we could in Exercises 5-4, No. 12, that it is global. As we already remarked in Section 4-8, the problem may define not one but several implicit functions and it is possible that a global extremum does not exist, or is an endpoint maximum for an implicit function other than the one for which a local extremum was found. An example is given by the equation

$$x^3 - x^2y - x^2 - xy + y^2 + y = (y - x^2)(y - x + 1) = 0$$

which has as its graph the parabola $y = x^2$ and the line $y = x - 1$ (Figure 5-6d). Our technique locates a local minimum of y at $x = 0$ on the branch $y = x^2$ but misses the fact that the graph considered as a whole does not have a lowest point. In general, before hard and fast conclusions can be drawn for any given problem of this type a deeper investigation is necessary.

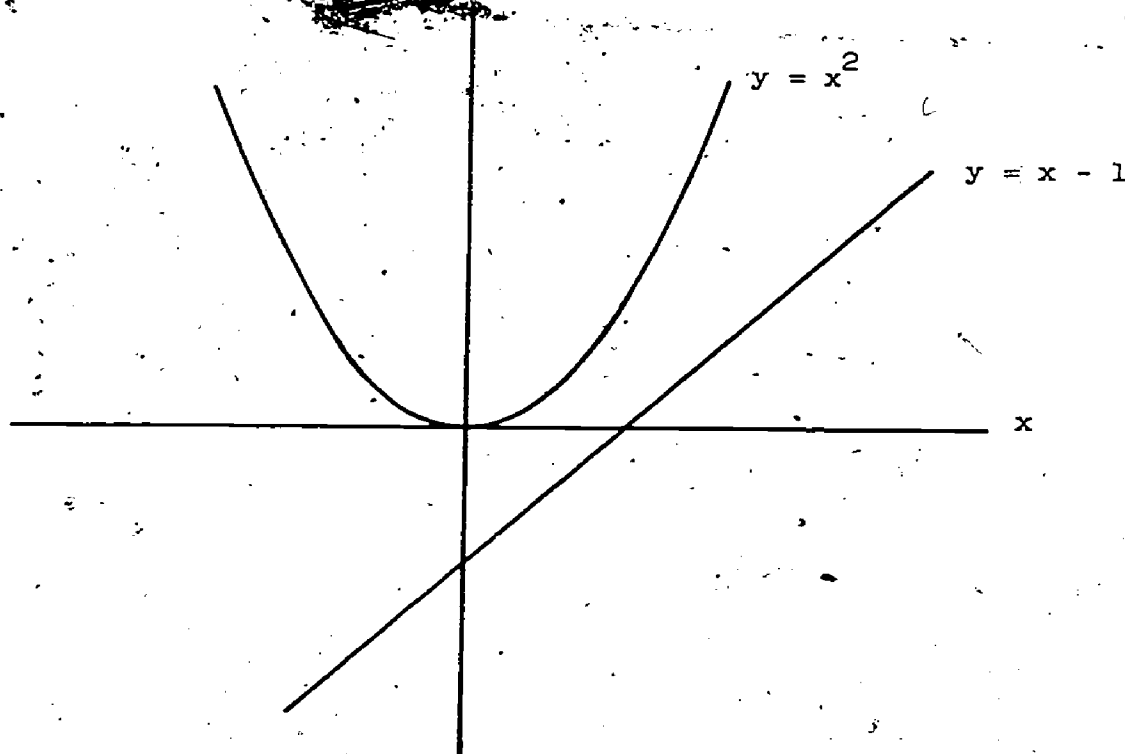
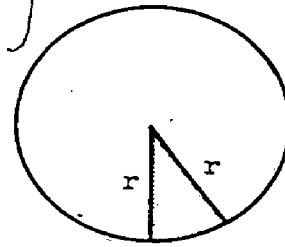


Figure 5-6d

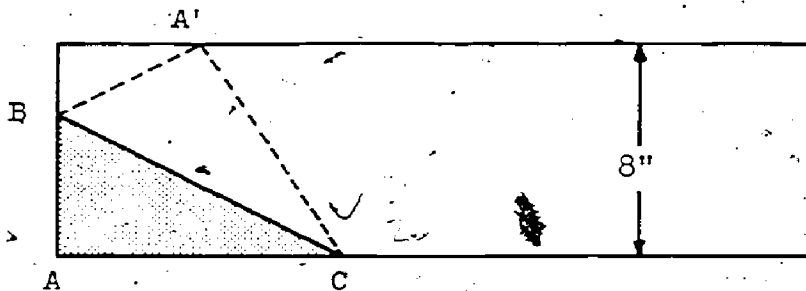
Exercises 5-6

1. (a) For a given volume V find the dimensions of a right cylindrical tin can with smallest surface area.
- (b) If the cost of the sides, top, and bottom of the can is a cents per square inch, and if the cost of the bead joining the top and bottom to the side is b cents per linear inch, find the most economical dimensions of the can for a given volume V .
2. (a) A cylindrical sheet-iron tank without top is to have volume V . Let h be the height of the tank, and r , the radius of the base. The side of the tank is to be constructed from stock costing P dollars per square foot and the base from stock costing Q dollars per square foot. Find the radius and height of the tank for which the cost of material in the tank is minimized.
- (b) More realistically, suppose that the base has to be cut from a square of side $2r$. Find the dimensions yielding minimum cost including the cost of material trimmed away.

3. Find the radius and height of the cone of greatest volume that can be made from a circular sheet of radius r by cutting a wedge from the center and bending the remaining portion to form a cone.



4. A point P is at a distance h above the center C of a sphere of radius r , where $h > r$. A cone is constructed having P for vertex, and for base the circle formed by cutting the sphere with a plane perpendicular to PC . In order to have the volume of the cone as great as possible, should this plane be above or below C ? How far?
5. A long strip of paper 8 inches wide is cut off square at one end. A corner A of this end is then folded over to the opposite side at A' , thus forming triangle ABC . Find the area of the smallest triangle that can be formed in this way.



6. (a) Let the general equation of a straight line L be given in the form $ax + by = c$.

Find the point Q on L for which the distance to a given point P not on L is a minimum. Prove that the line joining P to Q is perpendicular to L .

- (b) On the curve C given by $f(x,y)$ let Q be the point nearest to a point P not on the curve. If Q is not an endpoint of C , and all necessary derivatives exist, prove that the line joining P to Q is perpendicular to C .

7. Find the extrema of $x^2 + y^2$ if x and y are subject to the constraint $x^2 - 12x + y^2 - 8y + 51 = 0$. Give a geometric interpretation.
8. If $x_i \geq 0$, $i = 1, 2, \dots, n$ and $x_1 + x_2 + \dots + x_n = S$ (constant), find the maximum value of the product $x_1 \cdot x_2 \cdot \dots \cdot x_n$ (assuming it exists).

5-7. Tangent and Normal Lines.

In Section 5-3 we discussed the approximation of a differentiable function f on an interval by the linear function g which agrees with f at the end-points of the interval. We obtained estimates of the error with the aid of the Law of the Mean.

Here we consider a linear approximation to f in the neighborhood of a point a . For this purpose we use the linear function

$$h(x) = f(a) + f'(a)(x - a)$$

whose graph passes through $(a, f(a))$ and has the same direction as the graph of f at $x = a$. The line

$$y = f(a) + f'(a)(x - a)$$

is called the tangent to the curve $y = f(x)$ at $x = a$, (see Figure 5-7a).

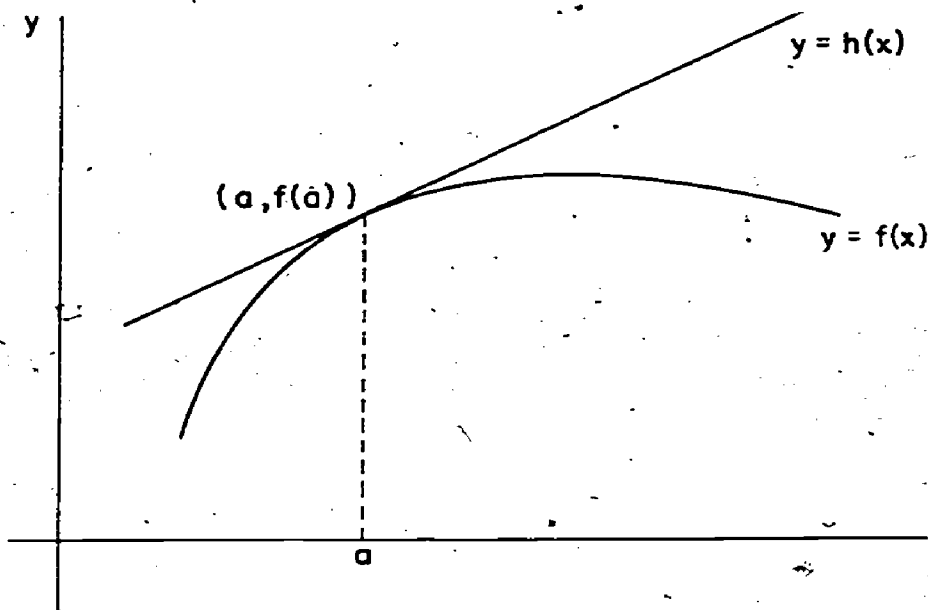


Figure 5-7a

Example 5-7a. The parabola $y = \frac{1}{2}x^2 + x + \frac{5}{2}$ goes through the point $(1, 4)$ with a slope $m = 2$. The line tangent there has the equation

$$y = 4 + 2(x - 1) \dots$$

Elimination of y from this equation and the equation of the parabola yields $\frac{1}{2}x^2 + x + \frac{5}{2} - (2x + 2) = \frac{1}{2}(x^2 - 2x + 1) = \frac{1}{2}(x - 1)^2 = 0$. It follows that the parabola and its tangent line meet only at $(1, 4)$.

Example 5-7b. We can always choose the coordinate axes so that the equation of a circle has the form $x^2 + y^2 = r^2$, where the positive constant r is the radius. Implicit differentiation gives, for the point (a, b) on the circle, the slope $m = -\frac{a}{b}$. The tangent line has the equation $y = b - \frac{a}{b}(x - a)$; we can rearrange algebraically, using the fact that $a^2 + b^2 = r^2$ to obtain the more symmetrical form

$$ax + by = r^2.$$

In geometry, the line tangent to a circle at one of its points was defined as the unique line through the point that meets no other point of the circle.* This definition would not serve for the parabola of Example 5-7a (Why not?). In a neighborhood of the point of tangency, the line tangent to a graph $y = f(x)$ does usually have the property that the curve lies on one side of the tangent. However, at special points (the points of inflection discussed in Section 5-8), the curve may cross its tangent. Thus the curve $y = x^3$ crosses its tangent at $x = 0$.

In case f has both first and second derivatives in the neighborhood of a we can easily obtain estimates for the error of approximation to the graph of f by the tangent line.

We find for the absolute error e , using the Law of the Mean,

$$\begin{aligned} e &= |f(x) - f(a) - f'(a)(x - a)| \\ &= |f'(z)(x - a) - f'(a)(x - a)| \end{aligned}$$

where z is some value between x and a . Thus,

$$e = |f'(z) - f'(a)| \cdot |x - a|,$$

and, applying the Law of the Mean again, we obtain

$$e = |f''(u)| \cdot |z - a| \cdot |x - a| \leq |f''(u)| \cdot (x - a)^2,$$

where u lies between z and a . If it is possible to obtain a bound for f'' in a neighborhood of a , say $|f''(x)| \leq M_2$, then we obtain the error estimate

$$(1) \quad e \leq M_2(x - a)^2.$$

This estimate for the error is to be compared with the estimate for the error of linear interpolation (Section 5-3,ii, Equation (6)).

* This property gave rise to the term tangent--derived from the Latin word tangens which means touching.

Example 5-7c. An estimate of the error in approximating $f(x) = x^3$ by the equation of the tangent at $x = a$ on the neighborhood $a - \delta \leq x \leq a + \delta$ is to be found.

We have $f''(x) = 6x$ for all x , so that

$$\begin{aligned} |f''(x)| &= 6|x| = 6|(x - a) + a| \\ &\leq 6(|x - a| + |a|) \\ &\leq 6(\delta + |a|). \end{aligned}$$

Consequently,

$$\begin{aligned} e &\leq 6(|a| + \delta)(x - a)^2 \\ &\leq 6(|a| + \delta)\delta^2. \end{aligned}$$

If, for example, we take $a = 5$, $\delta = .1$ we see that the error in estimating $f(x)$ is less than .31 in the approximation of a function value near $f(5) = 125$.

The tangent at a is the "best" linear approximation to the graph of f on the neighborhood of $x = a$ in a sense which is easily understood. If we take any other line passing through $(a, f(a))$, say

$$y = f(a) + m(x - a),$$

where $m \neq f'(a)$, then, from the Law of the Mean we obtain for the absolute error of approximation

$$\begin{aligned} \bar{e} &= |f(x) - f(a) - m(x - a)| \\ &= |f'(\bar{z}) - m| \cdot |x - a| \end{aligned}$$

where \bar{z} lies between x and a . Since f' is differentiable it is continuous, and by taking x sufficiently close to a we can make $f'(\bar{z})$ as close to $f'(a)$ as we wish. Since $m \neq f'(a)$ we can guarantee that $|f'(\bar{z}) - m|$ is greater than some fixed positive quantity K (say, $K = \frac{1}{2}|f'(a) - m|$) by taking $|x - a|$ small enough. Denote this quantity by K . We then have

$$(2) \quad \bar{e} > K|x - a|; \quad (x \neq a).$$

It follows from (1) and (2) that

$$e < \frac{M_2}{K}|x - a| \bar{e};$$

that is, by taking x sufficiently close to a , we can make the ratio of error e of approximation by the tangent to the error \bar{e} of approximation by any other line as small as desired.

In many problems we are concerned not only with the direction of the curve as represented by the tangent line, but also the direction perpendicular to it. For instance, each water particle at the front edge of a wave advancing on a beach moves along a path perpendicular to the edge (at least approximately). The line through $(a, f(a))$ perpendicular to the tangent line is defined to be the line normal to the curve $y = f(x)$ at $(a, f(a))$, or just the normal at the point $(a, f(a))$ (Figure 5-7b).

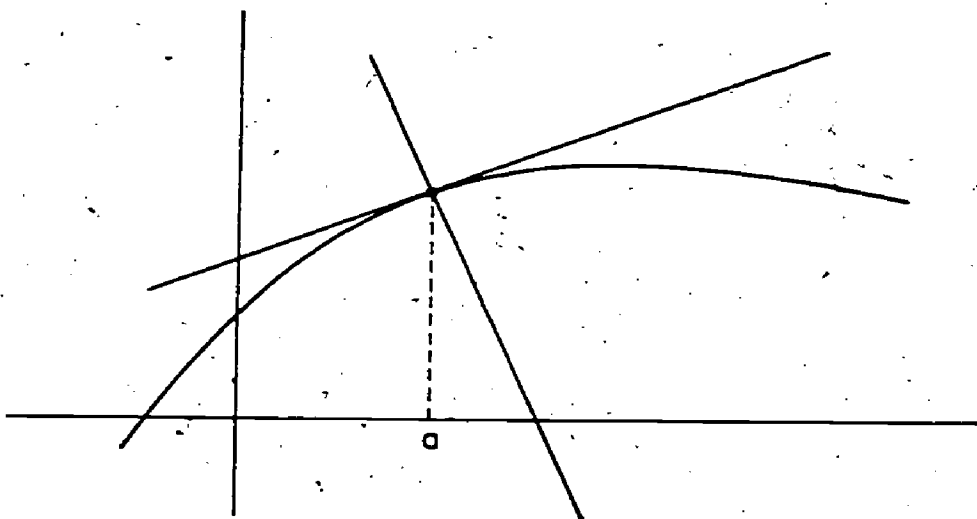


Figure 5-7b

If $f'(a) = 0$, then the normal is simply described by $x = a$. Otherwise, since the normal is perpendicular to the tangent, the slope $f'(a)$ of the curve and the slope m of the normal must satisfy the relation

$$(3) \quad mf'(a) = -1$$

and the normal has the equation

$$y = f(a) - \frac{(x - a)}{f'(a)} \quad *$$

The equation for the normal may also be written in the form

$$x = a - f'(a)[y - f(a)]$$

which is valid for all values of $f'(a)$, zero included.

*The relation (3) is proved as Theorem 2-3b, SMSG Intermediate Mathematics, p. 133.

Exercises 5-7

1. Show that the number of tangent lines that can be drawn from the point (h, k) to the curve $y = x^2$ is two, one, or zero according as k is less than h^2 , equal to h^2 , or greater than h^2 , respectively.

2. Find equations for the tangent and normal lines to the graphs of the following functions at the given points.

(a) $f : x \rightarrow x \sin x$, $x = 0$, $x = \frac{\pi}{2}$.

(b) $f : x \rightarrow \arcsin\left(\frac{1}{x}\right)$, $x = 2$.

(c) $f : x \rightarrow \frac{x^3}{1+x^2}$, $x = -1$, $x = 0$.

(d) $f : x \rightarrow \frac{x}{x^2 - 1}$, $x = x_0$ where $|x_0| \neq 1$.

3. (a) For the ellipse

$$\frac{x^2}{p^2} + \frac{y^2}{q^2} = 1,$$

and the hyperbola

$$\frac{x^2}{p^2} - \frac{y^2}{q^2} = 1,$$

obtain equations of the lines tangent at (a, b) on each curve in a symmetrical form like that of Example 5-7b for the circle.

- (b) For the same two curves, obtain equations for the normal lines at (a, b) on each curve in an analogous form.

4. Prove: The line tangent to the circle of Example 5-7b at (a, b) on the circle meets the circle at no other point.

5. Show that the graphs of the functions f and g where $f : x \rightarrow 6x^2$ and $g : x \rightarrow 4x^3 + 2$ have a common tangent line at the point $(1, 6)$. Sketch the graphs.

6. If $f : x \rightarrow ax^2 + bx + c$ ($a \neq 0$), show that the tangent line to the graph of f at the point $(p, f(p))$ is parallel to the chord joining the points $(m, f(m))$ and $(n, f(n))$ only if $p = \frac{m+n}{2}$.

7. Given the ellipse $b^2x^2 + a^2y^2 = a^2b^2$ and an arbitrary point P on the curve but not on either axis. Prove that if the normal at P to the ellipse passes through the origin, then the ellipse is a circle.

8. Comment on the footnote in Chapter 1, Section 1-1, page 6.

9. (a) Estimate the error of the approximation to $y = \sin x$ by the tangent at $x = 0$.

(b) Show that the error is at least third order in $|x|$; i.e., that

$$|x - \sin x| \leq C|x|^3$$

where C is constant.

10. Compare the methods of this section with the methods of Section 5-3 in Example 5-3b and Exercises 5-3, Numbers 15 and 16.

11. Show how to approximate $\sqrt[n]{2^n + 1}$ and estimate the error of approximation.

12. The location of an object at time t on a straight line is given by the law of motion

$$s = 5 \sin 3t - 3 \sin 5t$$

After it has started, when does the particle first reach a stop? How far is it then from the starting point?

13. Find an equation of the tangent line to the folium of Descartes

$$x^3 + y^3 - 3axy = 0$$

at the point (x_0, y_0) . Note particularly the situation at the point $(0, 0)$.

14. Find an equation of the tangent line to the graph of the equation

$$x^2 - x\sqrt{xy} - 2y^2 = 6$$

at the point $(4, 1)$.

5-8. Sketching of Graphs.

The problem we set for ourselves here is to obtain a sketch of the graph of a function, a picture which reveals the important general features of the graph but which need not be a precise point-by-point representation.

In coordinate geometry we saw that the x -intercepts (the solutions a of $f(a) = 0$ or zeros of f) and the y -intercept $f(0)$, yield easily plotted reference points $(a, 0)$ and $(0, f(0))$ on the graph of f (Int. Math. p. 145). Moreover, we learned tests for symmetry with respect to the y -axis and the origin, and observed that we can construct the entire graph of a symmetric function from the portion lying in a halfplane (Int. Math. pp. 147-148). We observed further (Section A2-1) that the graph of a function with period p can be constructed from the portion of the graph over any interval of length p .

From the calculus we obtain more information. Most of the functions we deal with here are differentiable to all orders on an interval. Such a function must be continuous and we know that its graph has no gaps. Furthermore, since the first derivative is continuous the graph is smooth; in particular, there can be no corners like that of the graph $y = |x|$ at the origin. Almost always the zeros of the derivative are isolated; that is, each zero has a neighborhood in which no other zero of the derivative appears. From Theorems 5-2a and 5-2b we then know that the graph of f is strongly monotone between successive zeros of the derivative. Furthermore, by observing the rise and fall of the values of the function at successive zeros of the derivative we can determine which are extrema. With this information we can obtain an excellent idea of the appearance of the graph.

We may wish also to incorporate information from a study of the second derivative. Thus, it is geometrically intuitive that if a curve is convex in a neighborhood of $x = a$ then it does not cross the tangent at $x = a$ (see Exercises 5-8, No. 11a). The curve $y = x^3$ does cross its tangent at $x = 0$, but as x increases there is a transition at $x = 0$ from downward flexure for $x < 0$ to upward flexure when $x > 0$ (see Exercises 5-8, No. 11b). At such a point of transition we must have $f''(a) = 0$ if the second derivative exists. These considerations suggest that such special points be singled out for consideration in a description of the gross properties of a function. In particular, we introduce the concept of point of inflection:

DEFINITION 5-8. If f'' is strongly monotone in the neighborhood of a , and $f''(a) = 0$ then $(a, f(a))$ is defined to be an inflection point of f .

The two possible cases are illustrated in Figures 5-8a and 5-8b. For the most part the zeros of the second derivative will be isolated and the graph will in general consist of convex arcs separated by points of inflection where the flexure reverses sense.

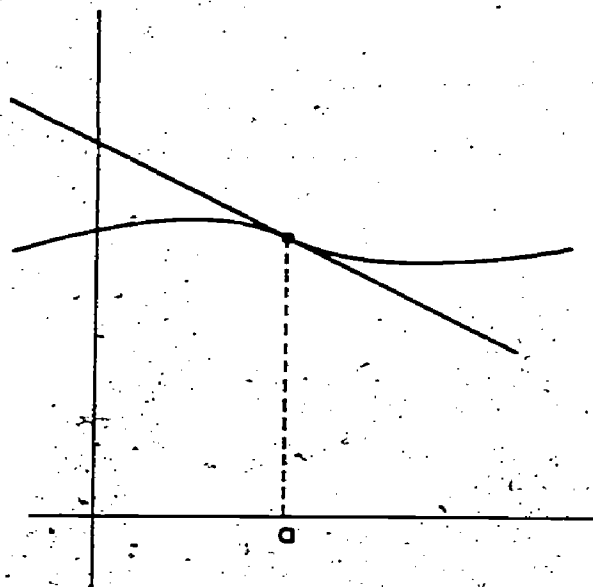


Figure 5-8a

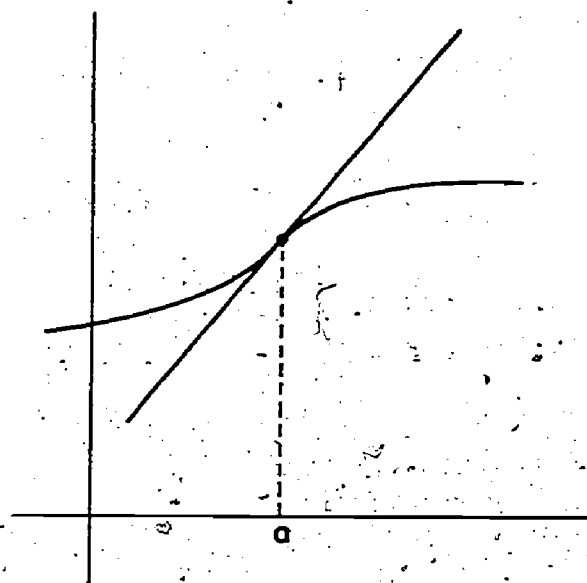


Figure 5-8b

Example 5-8a. The function $f(x) = 4x^5 + 5x^4 - 20x^3 - 50x^2 - 40x$ of Example 5-5b has the second derivative which may be factored:
 $f''(x) = 20(x + 1)^2(4x - 5)$. The graph of f is flexed downward for $x < -1$, and at $x = -1$ the second derivative is zero but the curve remains flexed downward as also may be seen from the fact that $f(-1)$ is a local maximum. At $x = \frac{5}{4}$, $f''(x)$ is again zero, and now it changes sign. There is an inflection point at $(1.25, -142.8)$, and the curve changes from downward to upward flexure, remaining flexed upward on the portion $x > \frac{5}{4}$ of its domain.

The concept of asymptote, a line approached by the graph for large values of x or y was introduced in the study of the hyperbola. We give simple criteria which help to locate horizontal and vertical asymptotes like those of the equilateral hyperbolas $xy = k$ (Int. Math. pp. 346-347).

In rough terms, an asymptote may be defined as a line which approximates a given curve, both in position and direction, at large distances from the origin. It is the idea of "large distances from the origin" which needs clarification.

The simplest case is that of a horizontal asymptote, $y = c$. The line $y = c$ may occur as an asymptote to the graph of x in two ways: for large positive values of x or for large negative values of x . In precise terms, $y = c$ is an asymptote of f for large positive values of x if for every $\epsilon > 0$ it is possible to make $|f(x) - c| < \epsilon$ by taking x large enough. In other terms, for each $\epsilon > 0$ we must find a value ω for which the condition is satisfied when $x > \omega$. For example, to prove that the hyperbola $xy = 1$ has the asymptote $y = 0$ for large positive values of x , we observe that

$$|f(x)| = \frac{1}{|x|} < \epsilon$$

provided $|x| > \omega \geq \frac{1}{\epsilon}$.

Clearly, in describing an asymptote we have defined a new kind of limit. We write

$$(1) \quad \lim_{x \rightarrow \infty} f(x) = c$$

if for every $\epsilon > 0$ there exists a value ω such that

$$|f(x) - c| < \epsilon$$

whenever

$$x > \omega.$$

The expression (1) is read "as x approaches infinity the limit of $f(x)$ is c ." The traditional word "infinity" does not signify anything mystical or vague here. It only means something in context, and in this case the context is precisely stated. In a similar vein, to describe asymptotes for large negative x , we say "as x approaches minus infinity the limit of $f(x)$ is c " and write

$$\lim_{x \rightarrow -\infty} f(x) = c$$

if for every $\epsilon > 0$ there exists a value ω such that

$$|f(x) - c| < \epsilon$$

whenever

$$x < \omega.$$

In this limit notation, the condition which defines $y = c$ as an asymptote to $y = f(x)$ for large positive x can be written

$$\lim_{x \rightarrow \infty} f(x) = c.$$

Example 5-8b. Let

$$f(x) = 4 + \frac{x}{(x+3)(x-1)}.$$

This is a typical case for a rational function. We find the horizontal asymptotes, if any, by comparing the leading terms of the numerator and denominator (see Exercises 5-8, No. 12). In the fraction above, the degree of the denominator is greater than that of the numerator and we conclude that the fraction approximates zero for large positive or negative x . It follows that f has the asymptote $y = 4$ for both large positive and negative x . The precise epsilon proof is left as an exercise.

A vertical asymptote $x = a$ can only occur at a point, a , where f is discontinuous because $|f(x)|$ exceeds any given positive real value for x in some sufficiently small deleted neighborhood of a . For a vertical asymptote, it is sufficient to show that $\frac{1}{f(x)}$ approximates zero when x is near a .

Example 5-8c. For the function f of the previous example we have

$$\frac{1}{f(x)} = \frac{(x+3)(x-1)}{(4x^2 + 9x - 12)},$$

and $\frac{1}{f(x)}$ approximates 0 for x near -3 and for x near 1 . The lines $x = -3$, $x = 1$ are vertical asymptotes (Figure 5-8c).

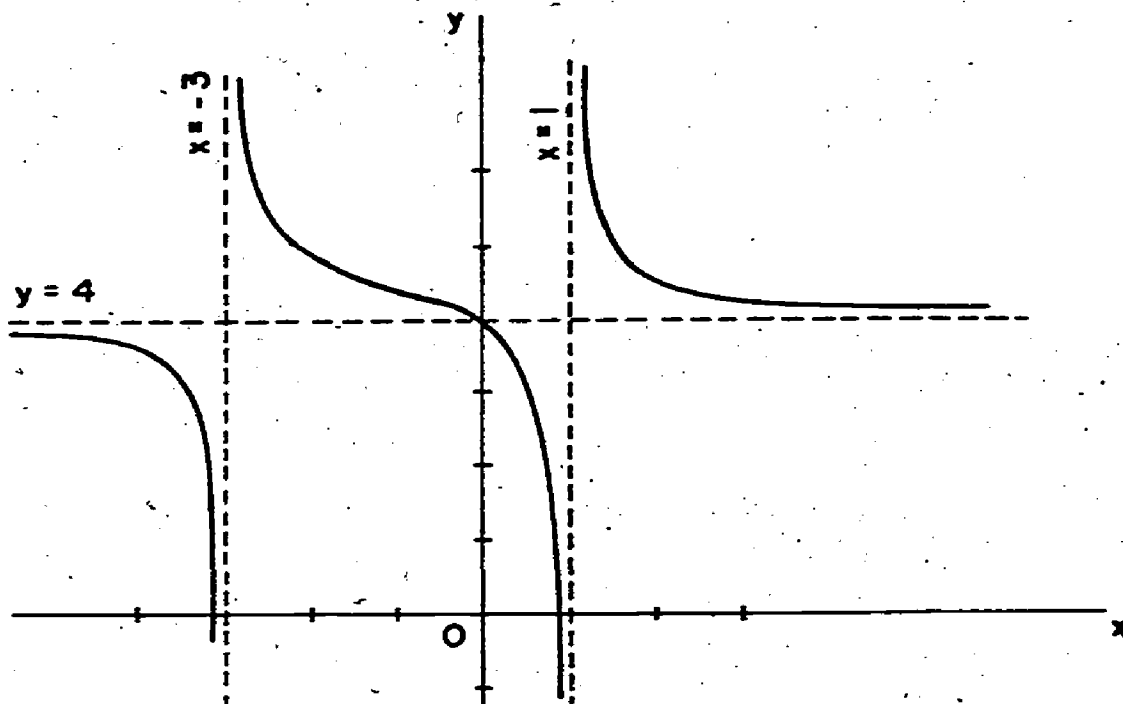


Figure 5-8c

Oblique asymptotes are also easily defined; the appropriate condition that the line $y = mx + b$ be an asymptote to the graph of f for large positive x is

$$\lim_{x \rightarrow \infty} |f(x) - [mx + b]| = 0.$$

If the graph of $y = f(x)$ has a slant asymptote it is easy to verify that the slope of the asymptote is given by

$$m = \lim_{x \rightarrow \infty} \frac{f(x)}{x}$$

(Exercises 5-8, No. 14a). On the other hand the limit may exist although the curve may not have an asymptote (Exercises 5-8, No. 14b).

With the wealth of auxiliary information obtained by the methods described above, the plotting of relatively few points is sufficient to obtain an adequate sketch of the graphs of most functions met in applications. Here we give a checklist of the information.

A. $f(x)$ determines:

1. the domain of definition.
2. the x-intercepts--the zeros of $f(x)$ --and the y-intercept-- $f(0)$.
(For a curve in implicit form $g(x,y) = 0$, the x-intercepts are the zeros of $g(x,0)$, and the y-intercepts are the zeros of $g(0,y)$.)

3. symmetry with respect to the y-axis if $f(x) = f(-x)$, with respect to the origin if $f(x) = -f(-x)$. (For a curve in implicit form there is symmetry with respect to the x-axis if $g(x, -y) = 0$ whenever $g(x, y) = 0$; with respect to the y-axis if $g(-x, y) = 0$ whenever $g(x, y) = 0$; with respect to the origin if $g(-x, -y) = 0$ whenever $g(x, y) = 0$.)

4. periodicity.

B. $f(x)$ also determines:

1. a horizontal asymptote $y = p$ to the right if

$$\lim_{x \rightarrow \infty} f(x) = p.$$

2. a horizontal asymptote $y = q$ to the left if

$$\lim_{x \rightarrow -\infty} f(x) = q.$$

3. a vertical asymptote $x = a$ if $\frac{1}{f(x)}$ approximates 0 in a deleted neighborhood of a .

4. a slant asymptote $y = mx + b$ to the right if both limits

$$m = \lim_{x \rightarrow \infty} \frac{f(x)}{x} \quad \text{and} \quad b = \lim_{x \rightarrow \infty} f(x) - mx$$

exist.

C. $f'(x)$ determines:

1. an interval on which $f(x)$ is increasing if $f'(x) > 0$ (weakly increasing if $f'(x) \geq 0$) on the interval.
2. an interval on which $f(x)$ is decreasing if $f'(x) < 0$ (weakly decreasing if $f'(x) \leq 0$) on the interval.
3. a maximum value $f(a)$ if $f'(a) = 0$ and $f'(x)$ changes from positive to negative at a .
4. a minimum value $f(a)$ if $f'(a) = 0$ and $f'(x)$ changes from negative to positive at a .

D. $f'(x)$ and $f''(x)$ determine:

1. a maximum value $f(a)$ if $f'(a) = 0$ and $f''(a) < 0$.
2. a minimum value $f(a)$ if $f'(a) = 0$ and $f''(a) > 0$.

E. $f''(x)$ determines:

1. an interval on which the graph of f is flexed downward if $f''(x) < 0$ on the interval.
2. an interval on which the graph of f is flexed upward if $f''(x) > 0$ on the interval.
3. an inflection point $(a, f(a))$ if $f''(a) = 0$ and $f''(x)$ changes sign as x increases through the value a .

It is a good procedure to prepare for the sketch of a graph by making a table in which pertinent items from the above check list are presented. Do all calculations separately from the table, so that only those data are shown which go directly into the sketch.

Example 5-8d. We give a complete checklist for the function

$$f(x) = 4 + \frac{x}{(x+3)(x-1)} \text{ of the previous two examples:}$$

x-intercepts: $x = -3.19, 0.94$

y-intercept: $y = 4$

horizontal asymptote: $y = 4$ to left and to right

vertical asymptotes: $x = -3$ and $x = 1$

intervals of decreasing $f(x)$: $x < -3$, $-3 < x < 1$, $1 < x$

downward flexure: $x < -3$, $x_1 < x < 1$, where x_1 is the abscissa of the inflection point

upward flexure: $-3 < x < x_1$, $1 < x$

inflection point: $(x_1, f(x_1))$, where $x_1^3 + 9x_1 + 6 = 0$; to one place after the decimal point, $-0.7 < x_1 < -0.6$

The sketch, Figure 5-8c, has already made use of all this information.

Note that intercepts, extrema, or inflection points may be only approximately determinable, with an accuracy dependent on your skill in approximately solving the appropriate equation $f(x) = 0$, $f'(x) = 0$, or $f''(x) = 0$.

Example 5-8e. Draw a graph of $f(x) = \cos x - 2 \sin x$, for $0 \leq x \leq 2\pi$. Here $f'(x) = -\sin x - 2 \cos x$; $f''(x) = -\cos x + 2 \sin x = -f(x)$. Zeros of $f(x)$ (and simultaneously of $f''(x)$) are numbers x for which $\tan x = \frac{1}{2}$; zeros of $f'(x)$ are numbers x for which $\tan x = -2$. Numerical values have been taken from a trigonometric table, rounded off to two places after the decimal point.

x-intercepts: $x = 0.46, x = 3.60$

y-intercept: $y = 1$

minimum: $x = 2.03, f(x) = -2.24$

maximum: $x = 5.17, f(x) = 2.24$

intervals of decreasing $f(x)$: $0 < x < 2.03, 5.17 < x < 2\pi$

intervals of increasing $f(x)$: $2.03 < x < 5.17$

inflection points: same as x-intercepts

intervals of downward flexure: $0 < x < 0.46$, $3.60 < x < 2\pi$

intervals of upward flexure: $0.46 < x < 3.60$

The points $f(x)$ for x a multiple of $\frac{\pi}{2}$ are also easily plotted, and were used in the construction of Figure 5-8d.

Finally, we sketch an implicitly defined curve (Figure 5-8d).

Note: the graph of $f : x \rightarrow \cos x - 2 \sin x$ could have been obtained more easily by noting that

$$\cos x - 2 \sin x = \sqrt{5} \left(\frac{\cos x}{\sqrt{5}} - \frac{2}{\sqrt{5}} \sin x \right) = \sqrt{5} \cos (x + \alpha) ,$$

where $\cos \alpha = \frac{1}{\sqrt{5}}$, but for the sake of illustration we have proceeded in a more complicated way.

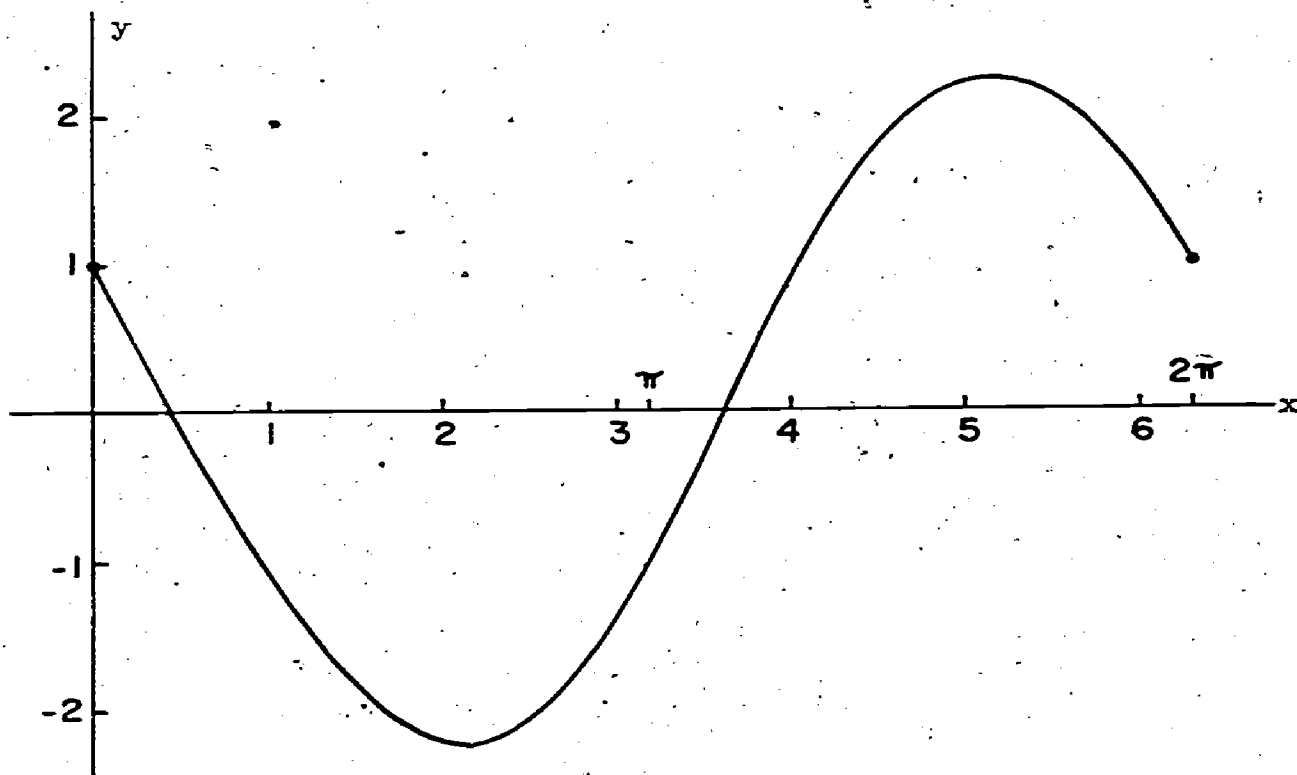


Figure 5-8d

Example 5-8f. Draw the curve with equation $\sqrt{x} + \sqrt{y} = 1$. Only points in the square $0 \leq x \leq 1$, $0 \leq y \leq 1$ can lie on the graph. Implicit differentiation gives

$$\frac{1}{2\sqrt{x}} + \frac{D_x y}{2\sqrt{y}} = 0$$

or

$$D_x y = -\sqrt{\frac{y}{x}},$$

and for the second derivative

$$\begin{aligned} D_x(D_x y) &= -\frac{1}{2} \sqrt{\frac{x}{y}} \left[\frac{x D_x y - y}{x^2} \right] = \frac{1}{2} \sqrt{\frac{x}{y}} \frac{\sqrt{xy} + y}{x^2} \\ &= \frac{1}{2} \frac{x + \sqrt{xy}}{x^2} = \frac{\sqrt{x} + \sqrt{y}}{2 x^{3/2}} = \frac{1}{2 x^{3/2}}. \end{aligned}$$

The checklist is quite short:

x-intercept: $y = 1$

y-intercept: $x = 1$

interval where y decreases: $0 < x < 1$

interval of upward flexure: $0 < x < 1$.

With the plotting of one additional point $(\frac{1}{4}, \frac{1}{4})$, the sketch is easily made (Figure 5-8f).

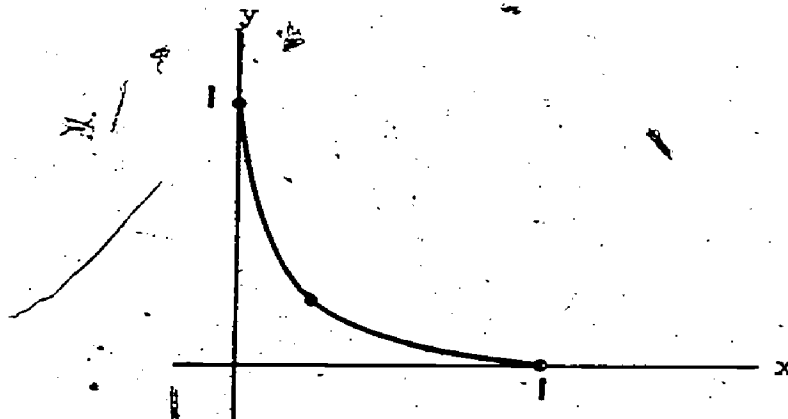


Figure 5-8e

Exercises 5-8

1. Draw the graph of

$$f : x \longrightarrow 4x^5 + 5x^4 - 20x^3 - 50x^2 - 40x.$$

(See Example 5-8a.)

2. Locate the point of inflection on the graph of $f : x \longrightarrow (x + 1) \arctan x$.

3. Determine equations of the horizontal and vertical asymptotes, if any, of the graph of

(a) $xy + y = x$.

(b) $x^2y - 3x + 2 = 0$.

4. Find horizontal and vertical asymptotes, maxima, minima, and inflection points: Show all tests used to identify each such point and draw the graph of the function.

(a) $f : x \longrightarrow \frac{x}{x^2 + 1}$

(b) $f : x \longrightarrow \frac{x^2}{x^2 + 1}$

5. Use information about extrema, and flexure to draw the graph of

$$f(x) = -x\sqrt{3 - x^2}.$$

6. Draw the graphs of the given functions making use of extrema and information about flexure:

(a) $f : x \longrightarrow \sqrt{\frac{x^3}{2 - x}}$.

(b) $f : x \longrightarrow \cos^2 x + 2 \cos x$.

(c) $f : x \longrightarrow x \arcsin x$.

7. Discuss symmetry, intercepts, asymptotes, extrema, intervals of flexure, and sketch the graph.

(a) $f(x) = \frac{2(x - 2)}{x^2}$

(b) $f(x) = x^{2/3}(x - 2)^2$

8. Draw the graph of $x^2y + xy^2 = 1$.
- Determine horizontal and vertical asymptotes, if any.
 - Locate the axis of symmetry and the point of intersection of the curve with this line.
 - Show that the curve has one and only one extremum and locate and classify this point.
9. Show that the function $x \mapsto \frac{ax + b}{cx + d}$ (assumed non-constant) has no maxima or minima regardless of the values of a , b , c , and d .
10. Prove that the inflection points of the graph of $y = x \sin x$ lie on the curve $y^2(4 + x^2) = 4x^2$.
11. (a) Prove that if a curve is differentiable and flexed downward, in an interval, the curve lies wholly under its tangent lines in this interval.
- (b) Show that the graph of f crosses its tangent at $x = a$ if a is a point of inflection.
- (c) Let $f(x) = x^n$ where n is a natural number. For what values of n , if any, does the graph have inflection points? Give sketches comparing the graphs for different values of n .
12. For a rational function given by
- $$r(x) = \frac{a_p x^p + a_{p-1} x^{p-1} + \dots + a_0}{b_q x^q + b_{q-1} x^{q-1} + \dots + b_0}$$
- where $a_p, b_q \neq 0$, find conditions for the existence of horizontal asymptotes, i.e., conditions that either limit, $\lim_{x \rightarrow \infty} r(x)$, $\lim_{x \rightarrow -\infty} r(x)$ exists.
13. For what values of a does the function $f : x \mapsto \frac{x^2 + 2x + a}{x^2 + 4x + 3a}$ assume all real values? Sketch the graph of the function for this case.

14. (a) Sketch the graph of the function

$$f: x \mapsto \frac{x^2}{x-a}, \quad a \neq 0,$$

and determine all the horizontal, vertical and slant asymptotes.

- (b) In part (a), the slope of the slant asymptote is $\lim_{x \rightarrow a} \frac{\frac{x^2}{x-a} - x}{x} = 1$.

Show that for the function

$$f: x \mapsto 1 + x - 2\sqrt{x}, \quad x \geq 0,$$

the $\lim_{x \rightarrow a} \frac{f(x)}{x}$ exists, although the graph has no straight line asymptotes.

Miscellaneous Exercises

1. Show that two tangent lines to the graph of $y = x^3 - 3x^2 + 3x$ pass through the point $(4, 1)$. Find their equations.

2. Show that the tangent line to the conic section

$$ax^2 + 2bxy + cy^2 + 2dx + 2ey + f = 0$$

at a point (x_0, y_0) on the curve has the equation:

$$ax_0x + b(y_0x + x_0y) + cy_0y + d(x_0 + x) + e(y_0 + y) + f = 0.$$

3. For what points (h, k) can one draw

- (i) two tangent lines,
- (ii) one tangent line,
- (iii) no tangent line

to the graph of

(a) $x^2 + 3xy + y^2 = a^2, \quad (a > 0).$

(b) $3x^2 + xy + 3y^2 = a^2, \quad (a > 0).$

(c) $\sqrt{x} + 5y = \sqrt{a}, \quad (a > 0).$

4. Determine equations of the horizontal and vertical asymptotes, if any, of the graph of

(a) $xy^2 - 4y - x = 0$.

(b) $xy - \cos x = 0$.

5. Sketch the graph of

(a) $y = \frac{(x-a)^m}{(x-b)^n}$; m, n integers, $m, n \geq 1$, $a \neq b$.

(b) $y^2 = \frac{(x-a)^m}{(x-b)^n}$; m, n integers, $m, n \geq 1$, $a \neq b$.

(c) $y = \frac{(x-a)^m(x-b)^n}{(x-c)^p}$; m, n, p integers, $m, n, p \geq 1$, $a \neq b \neq c$.

(d) $y^2 = \frac{(x-a)^m(x-b)^n}{(x-c)^p}$; m, n, p integers, $m, n, p \geq 1$, $a \neq b \neq c$.

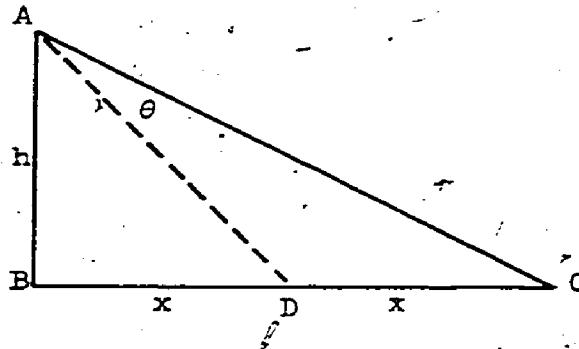
6. In past SMSG writers' conferences it has been observed that when the number of writers on a team is 28 or greater, all the available time is spent in discussion between members of the group, so that no writing gets done. Assuming that in a group of x writers ($28 \geq x \geq 1$) each participant is engaged in writing $40[1 - (\frac{x-1}{27})^2]$ hours per week, determine the size of the team which maximizes the total number of hours the group spent writing. (Draw your own conclusions about the team which wrote this book).

7. A picture h feet high is placed on a wall with its base b feet above the level of the observer's eye. If he stands x feet from the wall, verify that the angle of vision ϕ subtended by the picture is given by

$$\phi = \operatorname{arccot} \frac{x}{bh+b} - \operatorname{arccot} \frac{x}{b}.$$

Show that to get the "best" view of the picture, i.e., the largest angle of vision, the observer should stand $\sqrt{b(h+b)}$ feet away from the wall.

8. Let ABC be a right triangle with AB perpendicular to BC , the length of $AB = h$, the length of $BC = 2x$. Let AD be the median to side BC . Determine x so that the angle θ between the median and the hypotenuse of triangle ABC is a maximum.



9. The location of an object on a straight line at time t is given by the formula

$$S = At - (1 + A^4)t^2.$$

Show that the object moves forward initially, when A is positive, but ultimately retreats. Show also that for different values of A the maximum possible distance that the particle can move forward is $\frac{1}{8}$.

10. A man standing at the edge of a circular swimming pool wishes to reach a point $\frac{1}{4}$ of the way around the pool in the least possible time. He plans to run along the edge of the pool for some distance and then swim straight to his destination. If he can swim 20 feet per second and run 24 feet per second, how far should he run before diving in?
11. A conical cup with radius r , height h , is filled with water. Find the radius R of the sphere which displaces the largest volume of water when jammed into the cup.

Appendices

These appendices contain primarily three kinds of materials: matters which you are presumed to have studied before and may use for review or to acquire the necessary knowledge; matters alluded to in the text and amplified here; matters which are normally studied in later courses but supplied here for reference if you would like a brief account for completeness.

Appendix 1 contains reference material on the real number system. Appendix 2 contains reference material on functions, relations and their graphical representations. Appendix 3 describes the method of mathematical induction which is essential for mathematical literacy but is often neglected in the secondary curriculum. It is assumed that you are conversant with most of this material or will become so as the course progresses. The later appendices contain material which would be out of place in the text either because it constitutes an interesting digression or because it is normally studied at the next higher level. With the appendices many of the gaps left in the text may be reduced or eliminated altogether.

Appendix O

THE GREEK ALPHABET

A	α	alpha	N	ν	nu
B	β	beta	ξ	ξ	xi
Γ	γ	gamma	O	\omicron	omicron
Δ	δ	delta	Π	π	pi
E	ϵ	epsilon	P	ρ	rho
Z	ζ	zeta	Σ	σ	sigma
H	η	eta	T	τ	tau
Θ	θ	theta	Υ	υ	upsilon
I	ι	iota	Φ	ϕ	phi
K	κ	kappa	X	χ	chi
Λ	λ	lambda	Ψ	ψ	psi
M	μ	mu	Ω	ω	omega

Appendix 1

THE REAL NUMBERS

We assume that you are familiar with the real numbers represented as infinite decimals and as points of the number line (Intermediate Mathematics).

For the calculus, it is especially important that you master the material on inequalities and absolute value. If you are at all unsure of yourself, work through this material thoroughly. Many of the results obtained in the exercises are alluded to in the text.

A1-1 Algebraic Properties of the Real Numbers

We let \mathcal{R} denote the set of real numbers and represent the elements of \mathcal{R} by lower case (Roman) letters, a, b, c, \dots . The real numbers constitute an algebraic system in which addition, multiplication, and their inverse operations, subtraction and division (except by zero) can be performed. Such a system is known as a field. The field properties of \mathcal{R} are summarized in the following laws.

Closure

$$(1) \quad (a + b) \in \mathcal{R}, \quad ab \in \mathcal{R}.$$

If a and b are real numbers, then their sum and product are real numbers.

Commutative Laws

$$(2) \quad a + b = b + a, \quad ab = ba.$$

A sum or product of two real numbers is independent of the order in which they are taken.

Associative Laws

$$(3) \quad a + (b + c) = (a + b) + c, \quad a(bc) = (ab)c.$$

A sum or product of three real numbers is independent of the way in which they are associated in pairs.

Distributive Law

$$(4) \quad a(b + c) = ab + ac .$$

Multiplication is distributed over a sum.

Identities

There exist real numbers 0 (the additive identity) and 1 (the multiplicative identity) such that

$$(5) \quad a + 0 = a , a \cdot 1 = a .$$

Inverses

Each real number a has an additive inverse, $-a$, "the negative of a ", such that

$$(6a) \quad a + (-a) = 0 .$$

Each real number except 0 has a multiplicative inverse $\frac{1}{a}$, "the reciprocal of a ", such that

$$(6b) \quad a\left(\frac{1}{a}\right) = 1 , (a \neq 0) .$$

Finally, we require that 0 and 1 are not the same number.

$$(6c) \quad 0 \neq 1 .$$

Subtraction and Division

The operation of subtraction is defined by

$$(7a) \quad a - b = a + (-b) ,$$

and division by

$$(7b) \quad \frac{a}{b} = a\left(\frac{1}{b}\right) , (b \neq 0) .$$

We shall not attempt to derive the entire catalog of familiar properties of the real numbers implied by (1) - (7). The derivations of a number of these properties are left to you in exercises. There is one fact, however, to

which we wish to give prominent attention:

DIVISION BY ZERO CANNOT BE GIVEN ANY MEANING.

If $b \neq 0$, the statement (7b) assures us that $\frac{a}{b}$ is a real number; about the case $b = 0$ the statement is mysteriously silent. The fact is that it is impossible to assign a single definite value to $\frac{a}{b}$ consistent with (1) - (6). In proof we show first that

$$(8) \quad a \cdot 0 = 0.$$

We have

$$\begin{aligned} a \cdot 0 + a \cdot 0 &= a(0 + 0) && \text{(Distributive Law)} \\ &= a \cdot 0 && \text{(Additive Identity)} \end{aligned}$$

Put $\alpha = a \cdot 0$. By the preceding result

$$\alpha + \alpha = \alpha,$$

hence

$$(\alpha + \alpha) + (-\alpha) = \alpha + (-\alpha)$$

and

$$\alpha + (\alpha + (-\alpha)) = \alpha + (-\alpha) \quad \text{(Associative Law)}$$

and

$$\alpha + 0 = 0, \quad \text{(Additive Inverse)}$$

that is,

$$\alpha = 0. \quad \text{(Additive Identity)}$$

In order to define division by zero we must be able to find a "reciprocal of zero" under the definition of (6b), i.e., a number a such that

$$0 \cdot a = 1.$$

However, from (8) we have

$$0 \cdot a = 0$$

in contradiction to (6c). It follows that division by zero is meaningless and cannot be made meaningful. The attempt to define division by zero leads to a contradiction.

Exercises A1-1

1. Verify:
 - (a) that the natural numbers are not closed under subtraction and division.
 - (b) that the rational numbers form a field.
 - (c) that the operations of subtraction and division on the reals are not commutative or associative.
2. Prove: For any real number a , $-(-a) = a$.
3. Prove: For any real number a , $(-1)a = -a$.
4. Prove: For any real numbers a and b ,

$$-(a + b) = (-a) + (-b)$$
5. Prove: For any real numbers a and b ,

$$(-a)(-b) = ab$$
6. Prove: For any real numbers a and b ,

$$ab = 0 \text{ if and only if } a = 0 \text{ or } b = 0.$$

$$(\text{or } a = b = 0)$$
7. Verify that the numbers $\frac{a+b}{2}$, where a and b are rational numbers, constitute a field.

A1-2 Order Relations on the Real Numbers (Inequality)

The field relations are not only satisfied by the set \mathcal{R} of real numbers, but also by certain of its subsets including the set of rational numbers (see Exercise A1-1, No. 1b). Furthermore, \mathcal{R} itself is a subset of a number field, the field \mathcal{C} of complex numbers, which may be thought of as the set of all numbers of the form $a + bi$ where a and b are real and $i^2 = -1$. The real number system differs from the complex number system in one important respect: it is possible to order the reals by a relation of inequality. The properties of an order relation are given now in summary.

There is an order relation in \mathcal{R} , denoted by $a > b$ (read: "a is greater than b") with the following properties:

Trichotomy

Each a and b satisfy one and no more than one of the following relations:

$$a > b, \quad a = b, \quad b > a.$$

Transitive Law

If $a > b$ and $b > c$ then $a > c$.

Addition Law

If $a > b$, then $a + c > b + c$.

Multiplication Law

If $a > b$ and $c > 0$, then $ac > bc$; if $a > b$ and $c < 0$, then $bc > ac$.

It is often convenient to write $b < a$ (read "b is less than a") for $a > b$. If $a > 0$, a is said to be positive; if $a < 0$, a is said to be negative; thus, from the law of trichotomy a real number is positive, zero, or negative.

The two expressions $a > b$ and $b < a$ describe the same relation and neither is generally preferable to the other. We shall speak of $a > b$ and $b < a$ as strong inequalities.

Two other relationships which we shall use are the inequalities $a \geq b$ and $b \leq a$ (read "a is greater than or equal to b" and "b is less than or equal to a", respectively). The first of these, $a \geq b$, means that either $a > b$ or $a = b$; the second, $b \leq a$, means that either $b < a$ or $b = a$: the two inequalities represent exactly the same relation. (Sometimes, for emphasis, the relation $a \leq b$ is called a weak inequality.)

By the law of trichotomy, each a and b in \mathcal{R} satisfies one and no more than one of the following relations:

$$a \geq b, \quad b > a.$$

Thus, we note that if $a > b$, then $a \geq b$; also, if $a = b$, then $a \geq b$.

By the transitive property:

- (i) If $a \geq b$ and $b \geq c$, then $a \geq c$.
- (ii) If $a \geq b$ and $b > c$, then $a > c$.
- (iii) If $a > b$ and $b \geq c$, then $a > c$.

We observe that in (ii) and in (iii) the symbol " $>$ " representing strong inequality appears, hence we use the strong inequality in the conclusion " $a > c$ ".

In more general application of the transitive property the symbols " $>$ ", " \geq ", and " $=$ " may appear several times in a chain of reasonings; care must be exercised in selecting the symbol in the last step to insure that the conclusion is valid.

Example A1-2a:

If $a \geq b$, $b \geq c$, and $c = d$, then $a \geq d$.

It is convenient to write

$$a \geq b \geq c \geq d,$$

so that the valid conclusion " $a \geq d$ " is apparent. (Note: we did not write $a \geq b \geq c = d$. The valid conclusion is not " $a = d$ "; therefore, we shall avoid this form.)

Example A1-2b:

If $a \geq b$, $b \geq c$, $c = d$, $d \geq e$, and $e > f$, then $a > f$.

We write

$$a \geq b \geq c \geq d \geq e > f$$

which shows that the strong inequality " $a > f$ " is valid. We observe that the symbol " $>$ " for strong inequality appears at least once in the chain and hence we may use it in the final step. Note that within the chain we avoided writing $c = d$, but used $c \geq d$ as in the preceding example.

In our discussion of the transitive property we have used the symbols \geq , $>$, and $=$, but the statements also hold if the symbols \geq and $>$ are replaced by \leq and $<$, respectively. A word of caution is in order. We consider the expression

$$b < a > c$$

as completely meaningless (an example of what you should never write).

We leave the well-known properties of order for you to derive as exercises, but there is one property which we derive here as a useful example of such derivations:

THE SQUARE OF A NON-ZERO REAL NUMBER IS POSITIVE.

If $a > 0$, then from the multiplication law

$$a \cdot a > 0 \cdot a,$$

that is,

$$a^2 > 0.$$

If $a < 0$, then from the addition law,

$$a + (-a) < 0 + (-a)$$

and from the properties of the additive inverse and identity we obtain

$$0 < -a.$$

Thus $(-a)$ is positive and by the preceding argument $(-a)^2 > 0$. We know, in general (see Exercises A1-1, No. 5), that $(-a)(-b) = ab$. Setting $b = a$ in the last relation, we have $(-a)^2 = a^2$. It follows that $a^2 > 0$ when a is negative, and our argument is complete.

Exercises A1-2

1. Prove: For any real number a :
 - (a) if $a > 0$, then $0 > -a$.
 - (b) if $0 > a$, then $-a > 0$.
2. Prove: For any real numbers a, b, c, d if $a > b$ and $c > d$, then $a + c > b + d$.
3. Prove: For any real numbers a, b, c , if $a > b$ and $c < 0$, then $bc > ac$.
4. Prove:
 - (a) For any positive numbers a and b , $a > b$ if and only if $a^2 > b^2$.
 - (b) For any negative numbers a and b , $a > b$ if and only if $b^2 > a^2$.
5. Prove: For any real number a and any positive number b ,
 - (a) $a^2 \geq b$ if and only if $a \geq \sqrt{b}$ or $a \leq -\sqrt{b}$.
 - (b) $a^2 \leq b$ if and only if $-\sqrt{b} \leq a \leq \sqrt{b}$.
6. Prove: If $a > b > 0$ and $c > d > 0$, then $ac > bd$.
7. Prove: For any real numbers a and b if $ab > 0$, then either both $a > 0$ and $b > 0$ or both $a < 0$ and $b < 0$.
8. Prove: For any real number a
 - (a) if $a > 0$, then $\frac{1}{a} > 0$.
 - (b) if $a < 0$, then $\frac{1}{a} < 0$.
9. For $bd < 0$, show that $\frac{a}{b} < \frac{c}{d}$ if and only if $ad > bc$.
10. Show that for two positive numbers a and b , if $a > b$, then $\frac{1}{a} < \frac{1}{b}$.
11. Prove that the complex numbers form a field \mathbb{C} and that there can be no order relation on \mathbb{C} .
12. The field $\mathbb{Q}(\sqrt{2})$ of numbers of the form $a + b\sqrt{2}$ where a and b are rational numbers has the ordering relation $>$ because $\mathbb{Q}(\sqrt{2})$ is a subset of \mathbb{R} . Show that $\mathbb{Q}(\sqrt{2})$ is also ordered by the relation $>$ where

$$a + b\sqrt{2} > c + d\sqrt{2} \quad \text{means that}$$

$$a - b\sqrt{2} > c - d\sqrt{2}$$

13. Show that for all real numbers x and y

$$x^2 + xy + y^2 \geq 0.$$

14. (a) Prove that

$$(x + y)^2 \geq 4xy.$$

- (b) For positive numbers a and b , show that the arithmetic mean is not less than the geometric mean which is, in turn, greater than or equal to the harmonic mean:

$$\frac{a+b}{2} \geq \sqrt{ab} \geq \frac{2ab}{a+b}.$$

When does equality hold in this relation?

15. Find all values of x for which

$$ax^2 + 2bx + c \geq 0,$$

$$a \neq 0.$$

Discuss all possible cases.

16. Observe that

$$(a_1x + b_1)^2 + (a_2x + b_2)^2 + \dots + (a_nx + b_n)^2 \geq 0;$$

then use the solution of Number 15 to prove the Cauchy inequality

$$(a_1b_1 + a_2b_2 + \dots + a_nb_n)^2 \leq (a_1^2 + a_2^2 + \dots + a_n^2)(b_1^2 + b_2^2 + \dots + b_n^2),$$

with equality, if, and only if, $a_r = kb_r$ or $b_r = 0$, for $r = 1, 2, \dots, n$ and k , some constant.

17. If a_1, a_2, \dots, a_n are positive numbers, show that

$$\frac{a_1 + a_2 + \dots + a_n}{n} \geq \frac{n}{\frac{1}{a_1} + \frac{1}{a_2} + \dots + \frac{1}{a_n}}$$

$$(\text{Arithmetic Mean}) \geq (\text{Harmonic Mean})$$

This generalizes Exercise 14b.

18. Prove the general triangle inequality

$$\sqrt{x_1^2 + x_2^2 + \dots + x_n^2} + \sqrt{y_1^2 + y_2^2 + \dots + y_n^2} \geq \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2 + \dots + (x_n - y_n)^2}$$

A1-3 Absolute Value and Inequality

The absolute value of a real number a , written $|a|$, is defined by

$$|a| = \begin{cases} a, & \text{if } a \geq 0 \\ 0, & \text{if } a = 0 \\ -a, & \text{if } a < 0. \end{cases}$$

If we think of the real numbers in their representation on the number line, then $|a|$ is the distance between 0 and a (Figure A1-3). In general, for any real numbers a and b , the distance between a and b is

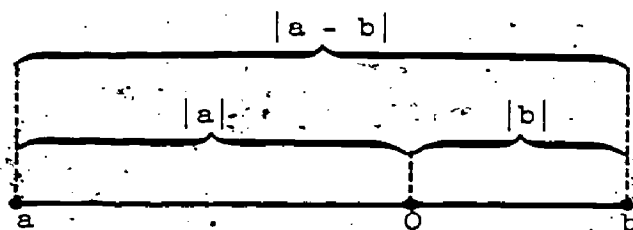


Figure A1-3

$|b - a| = |a - b|$. If x lies within the span $-\epsilon \leq x \leq \epsilon$ where $\epsilon \geq 0$, then clearly x is no farther from the origin than ϵ and we must have $|x| \leq \epsilon$. Conversely, if $|x| \leq \epsilon$, then $-\epsilon \leq x \leq \epsilon$. It follows immediately that

$$(1) \quad -|x| \leq x \leq |x|.$$

(See Exercises A1-3, No. 13a.)

From the inequalities

$$-|a| \leq a \leq |a| \quad \text{and} \quad -|b| \leq b \leq |b|$$

we obtain

$$-(|a| + |b|) \leq a + b \leq |a| + |b|,$$

whence

$$(2) \quad |a + b| \leq |a| + |b|.$$

(This relation is known as the "triangle inequality.") In words, the absolute value of a sum of two terms is not greater than the sum of the absolute value of the terms. Since any sum can be built up by successive additions, the result holds in general, viz.,

$$\begin{aligned} |a + b + c| &= |(a + b) + c| \\ &\leq |a + b| + |c| \\ &\leq |a| + |b| + |c|. \end{aligned}$$

We say that y is an upper estimate for x , and that x is a lower estimate for y if $x \leq y$. In (2) we have found an upper estimate for the absolute value of the sum $a + b$. It is often useful to have a lower estimate which is better than the obvious estimate 0. Such an estimate can be obtained from (2) by the device of setting $a = x + y$ and then setting $b = -x$ and $b = -y$ in turn. We then obtain

$$|y| - |x| \leq |x + y|$$

and

$$|x| - |y| \leq |x + y|.$$

Since $||x| - |y||$ is one or the other of the values $|x| - |y|$ or $|y| - |x|$, we have

$$(3) \quad ||x| - |y|| \leq |x + y|.$$

(See Exercises A1-3, No. 15.)

Special Symbols:

The symbol $\max\{r_1, r_2, \dots, r_n\}$ denotes the largest of the numbers r_1, r_2, \dots, r_n ; similarly, the symbol $\min\{r_1, r_2, \dots, r_n\}$ denotes the smallest of the numbers.

Example A1-3a:

$$\max\{2, 8, -3, -1\} = 8$$

$$\min\{2, 8, -3, -10\} = -10$$

$$\max\{-a, a\} = |a|.$$

Exercises A1-3

1. Find the absolute value of the following numbers:

(a) -1.75 .

(c) $\sin\left(\frac{-\pi}{4}\right)$.

(b) $\frac{-\pi}{4}$.

(d) $\cos\left(\frac{-\pi}{2}\right)$.

2. (a) For what real numbers x , does $\sqrt{x^2} = -x$?

(b) For what real numbers x does $|1 - x| = x - 1$?

3. Solve the equation:

(a) $|3 - x| = 1$.

(b) $|4x + 3| = 1$.

(c) $|x + 2| = x$.

(d) $|x + 1| = |x - 3|$.

(e) $|2x + 5| + |5x + 2| = 0$.

(f) $|2x + 3| = |5 - x|$.

(g) $2|3x + 4| + |x - 2| = 1 + |3 + x|$.

4. For what values of x is each of the following true? (Express your answer in terms of inequalities satisfied by x .)

(a) $|x| \leq 0$

(l) $|x - 1| + |x - 2| = 1$

(b) $|x| \neq x$

(m) $0 < |x^2 - a^2|$

(c) $|x| < 3$

(n) $|x - a| < \delta$

(d) $|x - 6| \leq 1$

(o) $0 < |x - a| < \delta$

(e) $|x - 3| > 2$

(p) $|x - 1| < 2$ and $|x + 1| < \frac{3}{2}$

(f) $|2x - 3| < 1$

(q) $|x - 1| < 2$ and $|2x - 1| < \frac{3}{2}$

(g) $|x - a| < a$

(r) $|x + y| = |x| + |y|$, for all y

(h) $|x^2 - 3| < 1$

(s) $|\sin x| = 0$

(i) $|(x - 2)(x - 3)| > 2$

(t) $|\sin x| > \frac{\sqrt{2}}{2}$

(j) $|x - 1| > |x - 3|$

(u) $|1 - \frac{1}{x}| < 1$

(k) $|x - 5| + 1 = |x + 5|$

(v) $\sqrt{|x|} > \frac{1}{2}$

5. Sketch the graphs of the following equations:

(a) $|x - 1| + |y| = 1$.

(b) $|x + y| + |x - y| = 2$.

(c) $y = |x - 1| + |x - 3|$.

(d) $y = |x - 1| + |x - 3| + 2|x - 4|$.

(e) $y = |x - 1| + |x - 3| + 2|x - 4| + 3|x - 5|$.

6. (a) Show that if $a > b > 0$, then

$$\frac{ab}{a+b} < b.$$

(b) Thus, show that for positive numbers a and b , the condition

$$\delta \leq \min\{a, b\} \text{ is satisfied by } \delta = \frac{ab}{a+b}$$

7. (a) Show for positive a, b that

$$\frac{a+b}{2} < \max\{a, b\} \text{ if } a \neq b.$$

(b) Prove for all a, b that

$$\max\{a, b\} = \frac{1}{2}(a + b + |a - b|).$$

(c) Prove for all a, b that

$$\min\{a, b\} = \frac{1}{2}(a + b - |a - b|).$$

8. Show that

$$\max\{a, b\} + \max\{c, d\} \geq \max\{a + c, b + d\}.$$

9. Show that if $ab \geq 0$, then

$$ab \geq \min\{a^2, b^2\}.$$

10. Show that if $a = \max\{a, b, c\}$, then $-a = \min\{-a, -b, -c\}$.

11. Denote $\min \left\{ \frac{a_1}{b_1}, \frac{a_2}{b_2}, \dots, \frac{a_n}{b_n} \right\}$ by $\min_r \left(\frac{a_r}{b_r} \right)$ and similarly for max.

If $b_r > 0$, $r = 1, 2, \dots, n$, prove that

$$\min_r \left(\frac{a_r}{b_r} \right) \leq \frac{a_1 + a_2 + \dots + a_n}{b_1 + b_2 + \dots + b_n} \leq \max_r \left(\frac{a_r}{b_r} \right).$$

12. Prove that

$$\frac{1}{n} \leq \frac{1 + 2 + \dots + n}{n^2 + (n-1)^2 + \dots + 2^2 + 1^2} \leq 1 \text{ for } n = 1, 2, 3, \dots$$

13. (a) Prove directly from the properties of order for $\epsilon > 0$ that if $-\epsilon \leq x \leq \epsilon$ then $|x| \leq \epsilon$. Conversely, if $|x| \leq \epsilon$ then $-\epsilon \leq x \leq \epsilon$.

(b) Prove that if x is an element of an ordered field and if $|x| < \epsilon$ for all positive values ϵ , then $x = 0$.

14. (a) Prove that $|ab| = |a||b|$.

(b) Prove that $\left|\frac{a}{b}\right| = \frac{|a|}{|b|}$, $b \neq 0$.

15. Prove that $|x - y| \leq |x| + |y|$.

16. Under what conditions do the equality signs hold for

$$||a| - |b|| \leq |a + b| \leq |a| + |b|?$$

17. If $0 < x < 1$, we can multiply both sides of the inequality $x < 1$ by x to obtain $x^2 < x$ (and, similarly, we can show that $x^3 < x^2$; $x^4 < x^3$, and so on). Use this result to show that if $0 < |x| < 1$, then $|x^2 + 2x| < 3|x|$.

18. Prove the following inequalities:

(a) $x + \frac{1}{x} \geq 2$, $x > 0$.

(b) $x + \frac{1}{x} \leq -2$, $x < 0$.

(c) $\left|x + \frac{1}{x}\right| \geq 2$, $x \neq 0$.

19. Prove: $x^2 \geq x|x|$ for all real x .

20. Show that if $|x - a| < \frac{|a|}{2}$, then

$$\frac{|a|}{2} < |x| < \frac{3|a|}{2}$$

for all $a \neq 0$.

21. Prove for positive a and b , where $a \neq b$, that

$$\frac{|b - a|^2}{4(a + b)} < \frac{a + b}{2} - \sqrt{ab} < \frac{|b - a|^2}{8\sqrt{ab}}$$

Al-4. Intervals, Neighborhoods

An interval of the number line is any set of real numbers which contains more than one member and which has the property that if a and b are elements of the set, then so is every real number x between them i.e., every x for which $a \leq x \leq b$. It is easy to verify that the following are intervals.

1. The segments (two endpoints). Given two real numbers a and b with $a < b$, there are four kinds of segments which have a and b as endpoints: the set of all x for which
 - (a) $a \leq x \leq b$, also written $[a, b]$ (closed interval)
 - (b) $a < x < b$, also written (a, b) (open interval)
 - (c) $a \leq x < b$, sometimes written $[a, b)$
 - (d) $a < x \leq b$, sometimes written $(a, b]$.

We emphasize, when the symbols $[a, b]$ and (a, b) are used for intervals it is assumed that $a < b$. (The sets defined by (c) and (d) are sometimes called half-open or half-closed intervals, but the concept is not particularly useful and we shall have no occasion to refer to it.)

2. The rays or half-lines (one endpoint). Given a real number a there are four kinds of rays having a as an endpoint: the set of all x for which
 - (a) $x \leq a$ or $x \geq a$, written $(-\infty, a]$ and $[a, \infty)$ (closed rays)
 - (b) $x < a$ or $x > a$, written $(-\infty, a)$ and (a, ∞) (open rays).

3. The entire number line, \mathbb{R} , $(-\infty, \infty)$.

A point of an interval which is not an endpoint is called an interior point of the interval.

The catalog above lists all types of intervals; it is possible (by use of the Separation Axiom, Section Al-5) to show that the catalog is exhaustive but we are not concerned with that question. In the text we distinguish only the closed and open intervals among the others; these are intervals having two endpoints and the endpoints are either simultaneously included in the interval or simultaneously excluded from the interval.

The length of an interval (whether open, closed, or other) with two endpoints a and b is the distance $|b - a|$ between a and b . The mid-point of an interval with endpoints a and b is the point $\frac{1}{2}(a + b)$. The closed interval with endpoints a and b is the set of values x satisfying

$$|x - \frac{1}{2}(a + b)| \leq \frac{1}{2} |b - a| ;$$

similarly, the open interval is the set of values satisfying the corresponding strong inequality

$$|x - \frac{1}{2}(a + b)| < \frac{1}{2} |b - a| .$$

Given $\delta > 0$, the δ -neighborhood of a real number a is the set of all points x satisfying

$$|x - a| < \delta$$

and δ is called the radius of the neighborhood. Thus the δ -neighborhood of a is the open interval

$$a - \delta < x < a + \delta$$

of length 2δ and midpoint at a . If we do not wish to specify the radius δ , we refer only to a neighborhood of a . Every open interval containing the point a contains a neighborhood of a ; conversely, every interval which contains with each of its points an entire neighborhood of the point, is open.

For many purposes it is useful to have the concept of a deleted neighborhood of a , that is, a δ -neighborhood with the center a deleted; namely, all x satisfying

$$0 < |x - a| < \delta .$$

Exercises A1-4

1. Use absolute value and inequalities to express the following facts.

- (a) The point x is closer to -2 than is the point a .
- (b) The point x is closer to point a than it is to the origin.

2. In each case, use absolute values and inequalities to express the fact that x is in the interval.
- (a) $[-5, 2]$
 - (b) $(-1, 5)$
 - (c) $[5.9, 6.1]$
 - (d) $(-2.95, -2.85)$
3. Find the interval or deleted interval to which all values of x must belong for each of the following:
- (a) $|x + 2| < 1$
 - (b) $0 < |x + 2| < 1$
 - (c) $|x + a| < \frac{|a|}{2}$
 - (d) $0 < |x + a| < \frac{|a|}{2}$
4. (a) A set of points is said to be bounded if there exists a real number A such that $|x| < A$, for all real members x of the set. Which of the intervals in Number 3 are bounded? Which are not? Prove your assertions.
- (b) A number M is said to be an upper bound for the set if $x \leq M$ for all members of the set; a number m is said to be a lower bound if $x \geq m$ for all members of the set. If a set has an upper bound, is it necessarily bounded? What if it has both an upper and a lower bound? If a set is bounded does it have both an upper and a lower bound? Prove your assertions.
5. For each of the following statements give the interval or intervals on which the statement is true.
- (a) $x^2 - x - 6 > 0$
 - (b) $(x - a)(x - b)(x - c) \leq 0$, for $a < b < c$.
 - (c) $\cos x > \sin x$
 - (d) $x + 1 \geq 2\sqrt{x}$
 - (e) $|x^2 - 1| < \frac{1}{100}$

6. In each of the following, for the given value of a find a neighborhood of a where the given inequality holds.

(a) $a = \frac{3}{2}$, $|2x - 3| < \frac{1}{7}$

(b) $a = \frac{\pi}{2}$, $|\sin x - 1| < 1 - \frac{1}{\sqrt{2}}$

(c) $a = -1$, $|x^2 + x| < \frac{1}{10}$

(d) $a = -1$, $|x^2 + x| < \frac{1}{100}$

(e) $a = -1$, $|x^2 + x| < \frac{1}{1000}$

A1-5. Completeness of the Real Number System. The Separation Axiom.

The field postulates and the postulates of order do not alone serve to define the real number system; the rational numbers satisfy the same postulates and so do other fields (Exercises A1-2, No. 12). Although no physical measurement requires anything more than the rational numbers, they are not adequate for either geometry or analysis. For example, the hypotenuse of a right triangle with legs of unit length has the irrational length $\sqrt{2}$; thus the Pythagorean Theorem would not exist if lengths were measured by rational values alone. In the rational field the concept of infinite decimal would be limited to terminating and periodic decimals; an infinite decimal like .101100111000 ... with chains of ones and zeros of increasing length is uninterpretable in the rational field. The system of rational numbers has theoretical gaps, but the real number system is complete in that real numbers are adequate to represent all the points on a line (lengths), and all infinite decimals. At the same time, it is possible to represent any real number by a point on a line or an infinite decimal; in fact, we use the concepts of point on the number line or infinite decimal as synonymous with real number.

The completeness of the real number system, its lack of theoretical gaps, is a consequence of a geometrically plausible axiom.

The Separation Axiom. If A and B are non-empty sets of real numbers for which every number in A is less than or equal to each number in B , then there is a real number s which separates A and B ; that is, for each $x \in A$ and $y \in B$ we have $x \leq s \leq y$.

In geometrical terms, if no point of a set A lies to the right of any point of a set B , then there is a point s such that all points of A (but s , should it happen to be a point of A) lie to the left of s , and all points of B (but s , if $s \in B$) lie to the right of s (see Figure A1-5a).

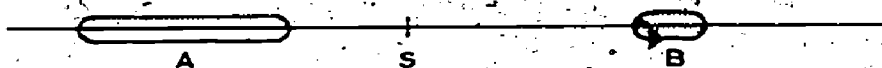


Figure A1-5a

A simple example of two sets satisfying the separation axiom is given by

$$A = \{x : x \leq -1\}, B = \{y : y \geq 1\}.$$

Clearly, any number s in the interval $[-1, 1]$ serves to separate these sets.

If two sets are separated by an entire interval, as in the preceding example, then it is possible to find a rational separation number s , because every interval on the number line contains rational points. The interesting cases are those for which there are elements of the two sets A and B closer together than any given positive distance. Gaps in the system of rational numbers can be exhibited as failures of the separation axiom for such sets. For example, let A be the set of positive rational numbers α satisfying $\alpha^2 < 2$, and let B be the set of positive rational numbers β satisfying $\beta^2 > 2$. It is possible to find rational values α and β closer together than any stated tolerance (see Exercises A3-1, No. 18) but a separation number s would have to satisfy $s^2 = 2$ and no rational number has that property (Exercise Al-5, No. 3c). We can define $\sqrt{2}$ as the unique real number which separates A and B . In fact, any real number can be defined as a separation number for suitable classes of rationals. More generally, it will be convenient for some purposes to determine a real number as the unique separation number for two sets by the criterion of the following lemma.

Lemma Al-5. Consider two sets of real numbers A and B such that $x \leq y$ for each $x \in A$ and each $y \in B$. If for every positive ϵ there exist $\alpha \in A$ and $\beta \in B$ such that $\beta - \alpha < \epsilon$, then the number s separating A and B is unique. Conversely, if there is just one separation number s , then for every positive ϵ there exist α and β with $\beta - \alpha < \epsilon$.

Proof. Let s and t be separation points for A and B . Given ϵ , $\alpha \in A$, and $\beta \in B$ such that $\beta - \alpha < \epsilon$, it follows from the fact that s and t lie between α and β (Figure Al-5b) that $|s - t| < \epsilon$. Since this

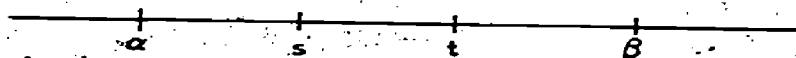


Figure Al-5b

is true for every positive ϵ it follows that $|s - t| = 0$ and hence that $s = t$ (see Exercises A1-3, No. 13b).

For the proof of the converse, let s denote the one number separating A and B . For every positive ϵ there must exist points $\alpha \in A$ and $\beta \in B$ such that

$$\alpha > s - \frac{\epsilon}{2} \quad \text{and} \quad \beta < s + \frac{\epsilon}{2},$$

for should one of these inequalities fail, then we would have $s - \frac{\epsilon}{2}$ or $s + \frac{\epsilon}{2}$ as a separation number. We conclude that $\beta - \alpha < \epsilon$.

Next we derive some important consequences of the Separation Axiom.

The Nested Interval Principle. A set of closed intervals $[a_n, b_n]$, $n = 1, 2, 3, \dots$, is called nested if $[a_{n+1}, b_{n+1}] \subset [a_n, b_n]$ for each natural number n . The principle states that there is at least one point s common to all the intervals of a nested set of closed intervals.

Proof. Let A consist of all the lower endpoints a_n of intervals of the nested set, and let B consist of all the upper endpoints b_n . The sets A and B satisfy the conditions of the Separation Axiom, and there exists a number s separating the two. Thus $a_n \leq s \leq b_n$.

The Least Upper Bound Principle. Let A be a set of numbers which is bounded above; i.e., there exists a value M such that $\alpha \leq M$ for all $\alpha \in A$. In the set of all upper bounds of A there is one upper bound which is smaller than any other, the least upper bound.*

Proof. Let B denote the set of upper bounds of A . The sets A and B satisfy the conditions of the Separation Axiom. It follows that there exists at least one separation number for A and B . Let s be such a separation number. Since s is a separation number it is an upper bound of A and is by definition an element of B . Since s is also a lower bound for B it is the least element of B and therefore the least upper bound of A .

The Least Upper Bound and Nested Interval Principles are also ways of expressing the completeness of the real numbers; they are equivalent to the Separation Axiom in the sense that either may replace the axiom and that the

*This number is also called the supremum of A and is denoted by $\sup A$. The abbreviation lub A is also common.

separation property will then follow.

In order to verify that the Separation Axiom and the Least Upper Bound Principle are equivalent formulations of the completeness of the real number system it is necessary to prove that in an ordered field the Least Upper Bound Principle implies the Separation Axiom. The proof is left as an exercise.

Corollary 1. If M is the least upper bound of the set A , then for each positive ϵ there exists an $\alpha \in A$ such that $\alpha > M - \epsilon$.

Corollary 2. A set of numbers which is bounded below has a greatest lower bound.

The proofs of these Corollaries are left as exercises.

Exercises A1-5

1. Prove Corollary 1 to the Least Upper Bound Principle.
2. Prove Corollary 2 to the Least Upper Bound Principle.
3. (a) Consider the sets A of positive rational numbers α satisfying $\alpha^2 < 2$, and B of positive rational numbers β satisfying $\beta^2 > 2$. Prove if $\alpha \in A$ and $\beta \in B$ that $\alpha < \beta$.
 (b) Show that a separation number s for the sets A and B must satisfy $s^2 = 2$; i.e., $s = \sqrt{2}$.
 (c) Prove that $\sqrt{2}$ is irrational.
4. (a) Prove for every real number a , that there is an integer n greater than a (Principle of Archimedes).
 (b) Prove that given any $\epsilon > 0$ there is an integer n such that $0 < \frac{1}{n} < \epsilon$.

5. (a) We define the infinite decimal

$$c_0.c_1c_2c_3 \dots,$$

where c_0 is an integer, and c_1, c_2, c_3, \dots , are digits, by the number r , where

$$c_0 + \frac{c_1}{10} + \frac{c_2}{10^2} + \dots + \frac{c_n}{10^n} \leq r < c_0 + \frac{c_1}{10} + \frac{c_2}{10^2} + \dots + \frac{c_n + 1}{10^n}.$$

Show that the preceding inequality does, in fact, define a unique real number.

(b) Given a real number r we define its decimal representation recursively in terms of the integer part function $[x]$ as follows:

$$c_0 = [r]$$

$$c_n = [10^n(r - c_0 - \frac{c_1}{10} - \frac{c_2}{10^2} - \dots - \frac{c_{n-1}}{10^{n-1}})].$$

Show that the inequality in part (a) is satisfied for this choice of c_n .

Show also that decimals consisting entirely of 9's from some point on are avoided. (Thus, we obtain $2 = 2.000 \dots$ but not $2 = 1.999 \dots$).

6. An infinite decimal $c_0.c_1c_2c_3 \dots$ is said to be periodic if for some fixed value p , the period of the decimal, we have $c_{n+p} = c_n$ for all n satisfying $n \geq n_0$, where we require that p is the smallest positive integer satisfying this condition. In words, from some place on, the decimal consists of the indefinite repetition of the same p digits. Thus

$$\frac{1}{3} = .33333\dots$$

$$\frac{15}{44} = .34090909\dots$$

are periodic decimals. It is convenient to indicate a cycle of p digits by underlining, rather than repetition; e.g.,

$$\frac{22}{7} = 3.\underline{142857} \dots$$

- (a) Prove that every periodic decimal represents a rational number.
(Hint: Consider the decimal as a geometric progression.)

- Λ(b) Prove that every rational number has a periodic decimal representation. (A "terminating" decimal in which each place beyond a certain point is zero is considered as a special case of periodic decimals.)
If $r = \frac{s}{t}$ represents a rational number given in lowest terms, find the largest possible period of the infinite decimal representation of r in terms of the denominator t .

From b we conclude that a decimal which is not periodic represents an irrational number, and conversely.

- Λ(c) Prove for every positive prime α other than 2 and 5 that there exists an integer, all of whose digits are ones, for which α is a factor; i.e., α is a factor of some number of the form

$$10^n + 10^{n-1} + 10^{n-2} + \dots + 10 + 1.$$

7. (a) Consider a polynomial with integer coefficients:

$$a_n x^n + a_{n-1} x^{n-1} + \dots + a_1 x + a_0. \quad (a_n \neq 0).$$

Prove that if $\frac{p}{q}$ is a rational root of this polynomial given in lowest terms, then p is a factor of a_0 and q is a factor of a_n .

- (b) Show that $x^3 + x + 1$ has no rational root.
(c) Prove that if \sqrt{n} is rational then it is integral.
(d) Prove that $\sqrt{3} - \sqrt{2}$ is irrational.

8. Prove that an ordered field in which the Nested Interval Principle holds also obeys the Separation Axiom.
9. Prove that an ordered field in which the Least Upper Bound Principle obeys the Separation Axiom.

Appendix 2

FUNCTIONS AND THEIR REPRESENTATIONS

A2-I. Functions.

The concept of a function is basic in the study of calculus; it is imperative that you have a clear understanding of the concept as well as related matters such as functional notation, operations on functions (composition, inversion), and special classes of functions such as monotone functions, polynomial functions, the absolute value function, the circular (trigonometric) functions, etc.

The precise definition of function can be formulated in many ways: as a set of ordered pairs (usually, ordered pairs of numbers), as an association or correspondence between two sets, etc. But no matter what definition we choose, for a function three things are required: a set called its domain, a set called its range, and a way of selecting a member of the range for each member of the domain.

Example A2-1a. The multiplication of integers by 2 defines a function. The domain of this function is the set of all integers; the range of the function is the set of all even integers.

We choose to define a function as an association between elements of two sets; thus the function of Example A2-1a associates with each integer its double.

DEFINITION A2-1. If with each element of a set A there is associated exactly one element of a set B , then this association is called a function from A to B . The set A is called the domain of the function, and the set C of all members of B assigned to members of A by the function is called the range of the function.

In what follows we shall be exclusively concerned with functions whose domains are subsets of real numbers and whose ranges are also subsets of real numbers. More complicated functions (like "vector valued functions") may be built from these.

The range C may be the whole set B , in which case the function is called an onto function, or it may be a proper subset of B . In any case, we generally take for B the whole set of real numbers, since a function is usually specified before its range is considered.

It is common practice to represent a function by the letter f (other letters such as F , g , h , ϕ , etc, will also be used): If x is an element of the domain of a function f , then $f(x)$ denotes the element of the range which f associates with x . (Read for $f(x)$ "the value of the function f at x ," or simply " f at x ," or " f of x .") An arrow is used to suggest the association of $f(x)$ with x :

$$f : x \longrightarrow f(x)$$

(read " f takes x into $f(x)$ "). This notation tells us nothing about the function f or the element x ; it is merely a symbolic description of the relation between x and $f(x)$.

Example A2-1b. Consider a function f defined as follows: f takes each number of the domain into its square. Thus, if 3 is an element of the domain, then f takes 3 into 9 , or f associates 9 with 3 . Concisely, $f(3) = 9$. In general, if x represents any number in the domain of f , then f takes x into x^2 :

$$f : x \longrightarrow x^2, \text{ or } f(x) = x^2.$$

The function is not adequately defined until we specify its domain. If the domain is the set of all integers $\{\dots, -2, -1, 0, 1, 2, \dots\}$, then the range is a subset of the nonnegative integers, $\{0, 1, 4, 9, 16, \dots\}$. If we choose the set of all real numbers as domain, then a different function is defined, even though the rule of association is the same; in this case the range of the function is the set of nonnegative real numbers.

Observe that a function from A to B is a one-way association; the reverse association from B to A is not necessarily a function. In Example A2-1b, $f(3) = 9$, and $f(-3) = 9$, while the reverse association would assign both 3 and -3 to 9 , violating the definition of a function.

It is often useful to think of a function as a mapping, and we say that a function maps each element of its domain upon one and only one element of its range. In this vein, $f : x \longrightarrow f(x)$ can be read, " f maps x upon $f(x)$;" $f(x)$ is called the image of x under the mapping, and x is called a

preimage of $f(x)$. This notion is illustrated in Figure A2-1a, where elements of the domain A and range B are represented by points and the mapping is suggested by arrows from the points of the domain to corresponding points of the range.

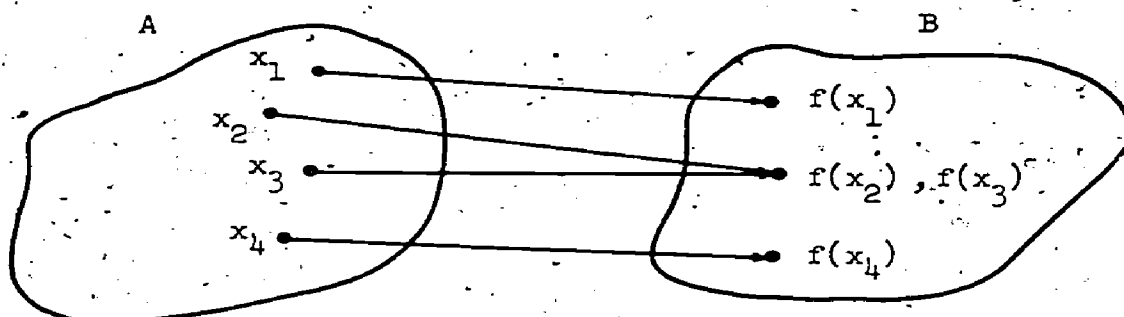


Figure A2-1a

Note that each element of the domain is mapped into a unique element of the range; i.e., each arrow starts from a different point in the domain. This is the requirement of Definition A2-1, that with each element of the domain there is associated exactly one element of the range.

Our Definition A2-1 of function contains the rather vague phrase, "there is associated." The manner of association must be specified whenever we are dealing with a particular function. In this course, a function will generally be defined by a formula giving its value: for example, $f(x) = 3x - 5$; $g(x) = x^2 + 3x + 7$. Other ways of defining a function include verbal description, graph, and table. (Recall Exercises 1-1.)

The notation $f(x)$ is particularly convenient when we refer to values of a function; i.e., elements in the range of the function. We illustrate this in the next example.

Example A2-1c. Consider the function

$$f : x \longrightarrow 3x^2 - 5$$

whose domain is the set of all real numbers. Then

$$f(x) = 3x^2 - 5,$$

$$f(-2) = 3(-2)^2 - 5 = 7,$$

$$f(0) = 3(0)^2 - 5 = -5,$$

and if $a + \sqrt{b}$ is a real number, then $f(a + \sqrt{b}) = 3(a + \sqrt{b})^2 - 5$.

We note, since x^2 may be any nonnegative real number, that $3x^2 - 5 \geq -5$, and hence the range of f is the set of all real numbers not less than -5 .

As mentioned earlier, a function is not completely defined unless the domain is specified. If no other information is given, it is a convenient practice, especially when dealing with a function defined by a formula, to assume that the domain includes all real numbers for which the formula describes a real number. For example, if a domain is not specified for the function $f : x \rightarrow \frac{2x}{x^2 - 9}$, then the domain is assumed to be the set of all real numbers except 3 and -3 . Similarly, if g is a function such that $g(x) = \sqrt{4 - x^2}$, we assume, in the absence of any other information, that the domain is $\{x : -2 \leq x \leq 2\}$; that is, the set of all real numbers x from -2 to 2 inclusive.

We note here that two functions f and g are identical if and only if they have the same domain and $f(x) = g(x)$ for each x in their domain.

The graph of a function is perhaps its most intuitively illuminating representation; it conveys important information about the function at a glance. The graph of f is the set of all those points (x, y) for which x is in the domain of f and $y = f(x)$.

Example A2-1d. The graph of the function $f : x \rightarrow y = \sqrt{25 - x^2}$ is the semicircle shown in Figure A2-1b.* The graph gives us a clear picture of what the function is doing to the elements of its domain, and we can, moreover, usually infer from the graph any limitations on the domain and range. Thus, it is easily determined from Figure A2-1b that the domain of f is the set of all x such that $-5 \leq x \leq 5$ and the range is the set of all y such that $0 \leq y \leq 5$. These sets are represented by the heavy segments on the x - and y -axis, respectively.

* In this figure a complete graph is displayed. The graph in Figure A2-1c, as well as most of the graphs in the text, are necessarily incomplete.

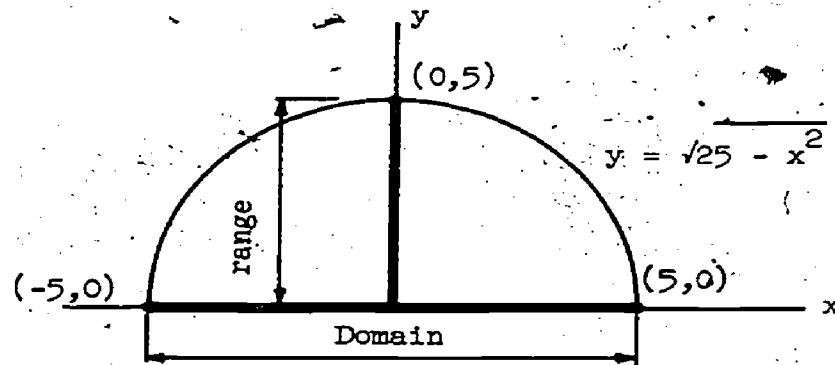


Figure A2-1b

We remind you of the fact that not every curve is the graph of a function. In particular, Definition A2-1 requires that a function map each element of its domain onto only one element of its range. In terms of points of a graph, this means that the graph of a function does not contain the points (x_1, y_1) and (x_1, y_2) if $y_1 \neq y_2$; i.e., two points having the same abscissa but different ordinates. This is the basis for the "vertical line test": if in the xy -plane we imagine all possible lines which are parallel to the y -axis, and if any of these lines cuts the graph in more than one point, then the graph represents a relation which is not a function. Conversely, if every line parallel to the y -axis intersects a graph in at most one point, then the graph is that of a function.

Example A2-1e. The equation $x^2 + y^2 = 25$, whose graph is a circle with radius 5 and center at the origin, does not define a function. On the open interval $-5 < x < 5$, every value of x is associated with two different values of y , contrary to the definition of function. Specifically, $(3, 4)$ and $(3, -4)$ are two points of the circle; they determine a line parallel to the y -axis and intersecting the circle in two points, thus illustrating that the circle is the graph of a relation that is not a function. We can, however, separate the circle into two semi-circles--the graphs of the functions

$$x \rightarrow \sqrt{25 - x^2} \quad (\text{Example A2-1d}) \quad \text{and} \quad x \rightarrow -\sqrt{25 - x^2}.$$

Throughout this discussion we have used the letters x and y to represent elements of sets. Specifically, if f is the function

$$f : x \rightarrow y = f(x),$$

then x represents an element (unspecified) in the domain of f , and y represents the corresponding element in the range of f . In many textbooks x and y are called variables, and since a particular value of y in the range depends upon a particular choice of x in the domain, x is called the independent variable and y the dependent variable. The functional relationship is then described by saying that " y is a function of x ." For the most part this language is not used in this textbook.

We conclude this section with a summary of several different special functions; you are undoubtedly acquainted with some of them.

The Constant Function. If b is an arbitrary real number, then the function f which associates with every real number x the value b , $f : x \rightarrow b$, is called a constant function. More generally, any function whose range contains exactly one number is a constant function. The graph of a constant function, say $f : x \rightarrow c$ for all real x , is a line parallel to and $|c|$ units from the x -axis.

The Identity Function. Let A be the set of all real numbers; with each number a in A , associate the number a . This association defines a function whose domain is A and whose range is A , namely

$$f : x \rightarrow x.$$

More generally, for any domain such a function is called the identity function. If the domain is the set of all real numbers, then the graph of f is the line with equation $y = x$.

The Absolute Value Function. With each real number the absolute value function associates its absolute value (Section A1-3):*

$$f : x \rightarrow |x| = \begin{cases} x & \text{for } x \geq 0, \\ -x & \text{for } x < 0. \end{cases}$$

* Alternative definitions:

$$f : x \rightarrow |x| = \max \{x, -x\};$$

$$f : x \rightarrow |x| = \sqrt{x^2}.$$

The graph of f is shown in Figure A2-1c; it is the union of two rays issuing from the origin.

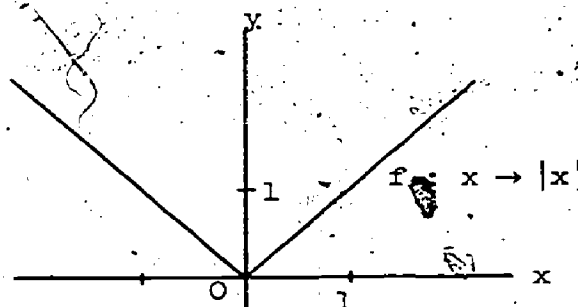


Figure A2-1c

The Integer Part Function. Every real number x can be represented as the sum of an integer n and a real number r such that

$$x = n + r, \quad n \leq x, \quad \text{and} \quad 0 \leq r < 1.$$

For example,

$$5.38 = 5 + .38,$$

$$3 = 3 + 0,$$

$$-2.4 = -3 + .6.$$

We call n the integer part of x and denote it by $[x] = n$; it follows that $[x] \leq x < [x] + 1$. Thus we see that to each real number x there corresponds a unique integer part $[x]$, and this correspondence defines the integer part function

$$f: x \longrightarrow [x].$$

A graph of this function is shown in Figure A2-1d; it is called a step graph; i.e., the graph of a step function.

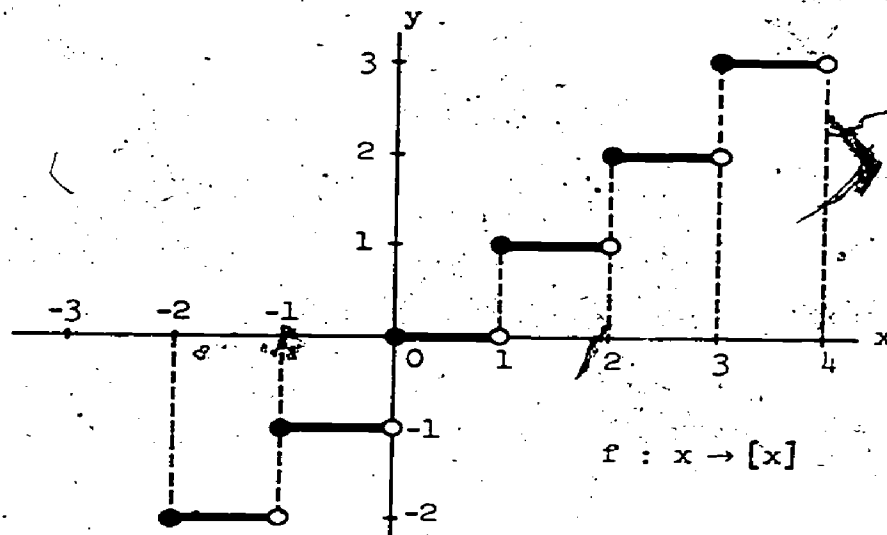


Figure A2-1d

The Signum Function. With each positive real number associate the number +1, with zero associate the number 0, and with each negative real number associate the number -1. These associations define the signum function, symbolized by $\text{sgn } x$. Thus

$$\text{sgn } x = \begin{cases} 1, & x > 0, \\ 0, & x = 0, \\ -1, & x < 0. \end{cases}$$

We leave it as an exercise for you to sketch the graph of this function.

Even and Odd Functions. Let f be a function whose domain contains $-x$ whenever it contains x . The function f is said to be even if $f(-x) = f(x)$. For example, the function f with values $f(x) = x^2$ is even since $(-x)^2 = x^2$ for all x . Geometrically the graph of an even function is symmetric with respect to the y -axis.

The function f is said to be odd if $f(-x) = -f(x)$. For example, the function f with values $f(x) = x^3$ is odd since $(-x)^3 = -x^3$ for all x . Geometrically the graph of an odd function is symmetric with respect to the origin.

Periodic Functions. Certain functions have the property that their function values repeat themselves in the same order at regular intervals over the domain (Figure A2-1e).

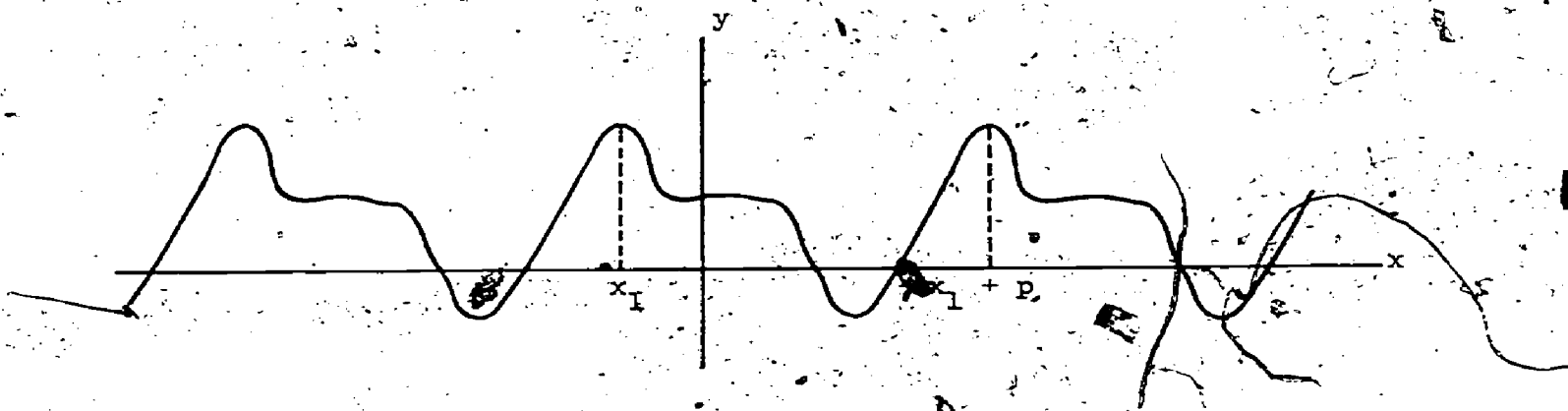


Figure A2-1e

Functions having this property are called periodic; included in this important class are the circular (trigonometric) functions, to be discussed in Section A2-5.

A function f is periodic and has period p , $p \neq 0$, if and only if, for all x in the domain of f , $x + p$ is also in the domain and

$$(1) \quad f(x + p) = f(x) .$$

From the definition we note that each successive addition or subtraction of p brings us back to $f(x)$ again. For example,

$$\begin{aligned} f(x + 2p) &= f((x + p) + p) \\ &= f(x + p) \\ &= f(x) , \end{aligned}$$

and

$$\begin{aligned} f(x - p) &= f((x - p) + p) \\ &= f(x) . \end{aligned}$$

In general, we infer that any multiple of a period of f is also a period; that is,

$$f(x + np) = f(x) \quad \text{for any integer } n .$$

For a constant function

$$f : x \rightarrow c ,$$

it is obvious that f is periodic with any period p , since

$$f(x + p) = c = f(x) .$$

It can be shown that for nonconstant periodic functions (continuous at one point at least)¹ there is a least positive value of p for which (1) is true. This is called the fundamental period, or simply the period, of such a function.

Example A2-1f. $f : x \rightarrow x - [x]$, x real, is a periodic function. If $x = n + r$ where n is the integer part of x and r its fractional part,² then

$$\begin{aligned} f(x) &= f(n + r) \\ &= (n + r) - [n + r] \\ &= n + r - n \\ &= r, \end{aligned}$$

and

$$\begin{aligned} f(x + 1) &= f(n + 1 + r) \\ &= (n + 1 + r) - [n + 1 + r] \\ &= n + 1 + r - (n + 1) \\ &= r. \end{aligned}$$

Thus, as was asserted, f is periodic and its period is 1, as shown in its graph (Figure A2-1f).

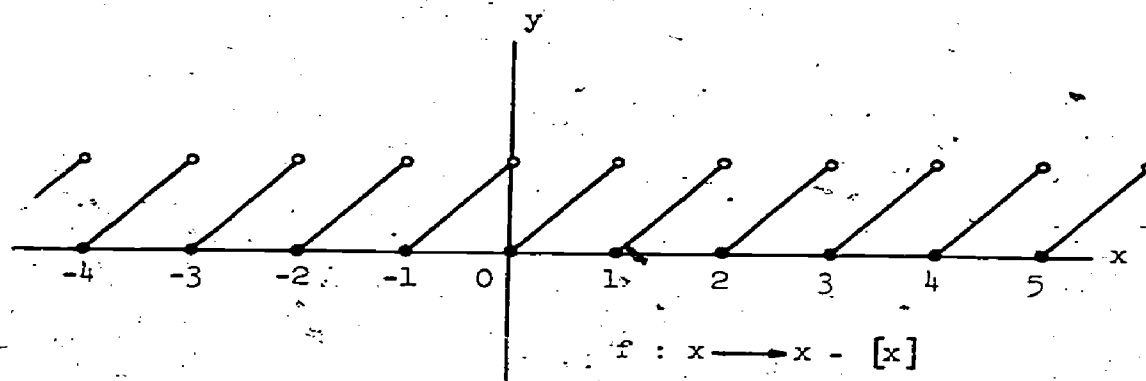


Figure A2-1f

¹This is left as Exercise 14(b) of 3-5 after the discussion of continuity.

²We note that since $f(x) = r$, the fractional part of x , this function is sometimes called the fractional part function.

Exercises A2-1

1. Below are given examples of associations between elements of two sets. Decide whether each example may properly represent a function. This also requires you to specify the domain and range for each function. Note that no particular variable has to be the domain variable, and also that some of the relations may give rise to several functions.

- (a) Assign to each nonnegative integer n the number $2n - 5$.
- (b) Assign to each real number x the number 7 .
- (c) Assign to the number 10 the real number y .
- (d) Assign to each pair of distinct points in the plane the distance between them.
- (e) $y = -3$ (for all x).
- (f) $x = -4$ (for all y and z).
- (g) $x + y = 2$.
- (h) $y = 2x^2 + 3$.
- (i) $y^2 - 4 = x$.
- (j) $y < 2x - 1$.
- (k) $f(x) = -\sqrt{16 - x^2}$.
- (l) $x^2 + y^2 = 16$.

2. Sketch the graphs of equations (e) - (l) of Number 1.
3. A function f is completely defined by the table:

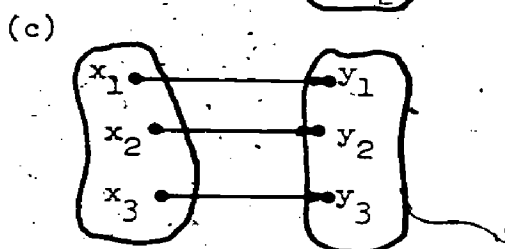
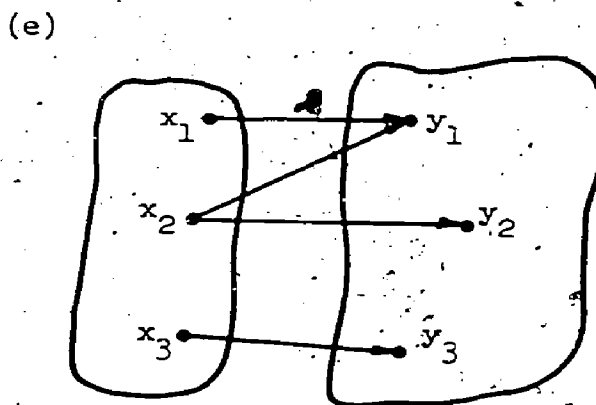
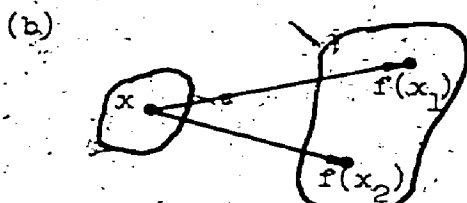
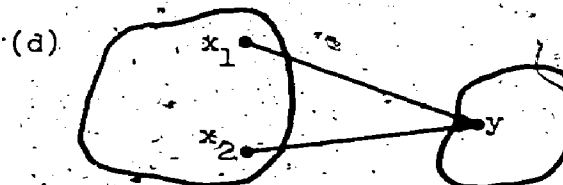
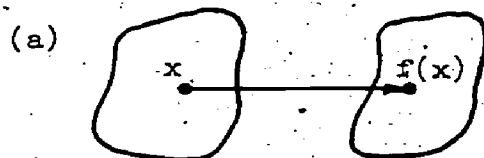
x	0	1	2	3	4
$f(x)$	-3	1	5	9	13

- (a) Describe the domain and range of f .
 - (b) Write an equation with suitably restricted domain that defines f .
4. If $f : x \mapsto x^2 + 3x - 4$, find
- (a) $f(0)$
 - (b) $f(2)$
 - (c) $f(-1)$
 - (d) $f(\sqrt{3})$
 - (e) $f(2 - \sqrt{2})$
 - (f) $f(f(1))$ (Hint: This is the value of f at $f(1)$.)

5. If g is a function defined by $g(x) = \frac{2x}{\sqrt{5-x^2}}$, find, if possible,

- (a) $g(0)$ (d) $g(2)$
 (b) $g(1)$ (e) $g(-3)$
 (c) $g(-1)$ (f) $g(\sqrt{5})$

6. Which of the following mappings represent functions?



7. Given the functions $f: x \rightarrow x$ and $g: x \rightarrow \frac{x^2}{x}$. If x is a real number, are f and g the same function? Why or why not?

8. Given the functions $f: x \rightarrow x+2$ and $g: x \rightarrow \frac{x^2-4}{x-2}$. If x is real, are f and g the same function? Why or why not?

9. What number or numbers have the image 10 under the following mappings?

- (a) $f: x \rightarrow 2x$ (d) $\alpha: x \rightarrow |x-4|$
 (b) $g: x \rightarrow x^2$ (e) $\phi: x \rightarrow [x]$
 (c) $h: x \rightarrow \sqrt{x^2+36}$

10. Which of the following statements are always true for any function f , assuming that x_1 and x_2 are in the domain of f ?

- (a) If $x_1 = x_2$, then $f(x_1) = f(x_2)$.
 (b) If $x_1 \neq x_2$, then $f(x_1) \neq f(x_2)$.
 (c) If $f(x_1) = f(x_2)$, then $x_1 = x_2$.
 (d) If $f(x_1) \neq f(x_2)$, then $x_1 \neq x_2$.

11. If $f(x) = |x|$, which of the following statements are true for all real numbers x and t ?

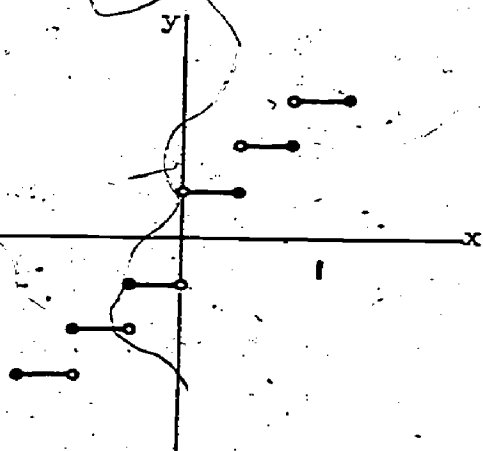
- (a) f is an odd function.
- (b) $f(x^2) = f(x)^2$.
- (c) $f(x - t) \leq f(x) - f(t)$.
- (d) $f(x + t) \leq f(x) + f(t)$.

12. Which of the following functions are even, which are odd, and which are neither even nor odd?

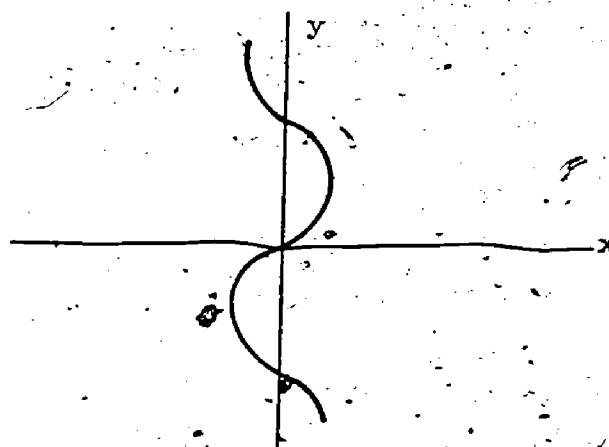
- (a) $f : x \rightarrow 3x$
- (b) $f : x \rightarrow -2x^2 + 5$
- (c) $f : x \rightarrow x^2 - 4x + 4$
- (d) $f : x \rightarrow -2x + 1$
- (e) $f : x \rightarrow x^3 + 4$
- (f) $f : x \rightarrow x^3 - 2x$
- (g) $f : x \rightarrow 2^{1/x}$
- (h) $f : x \rightarrow 2^{1/x^2}$

13. Which of the following graphs could represent functions?

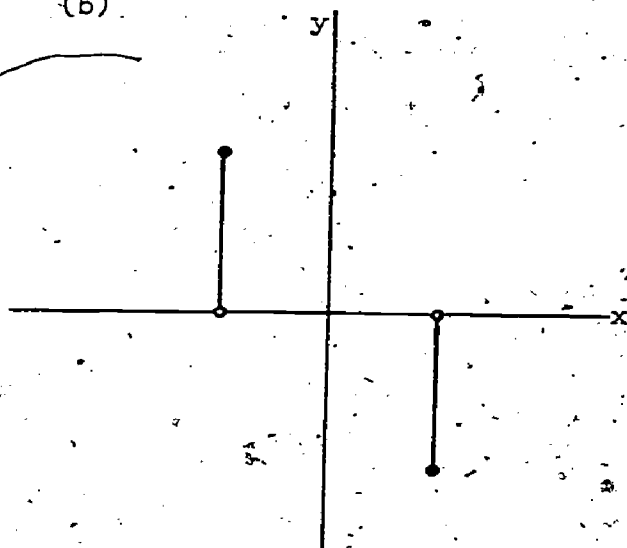
(a)



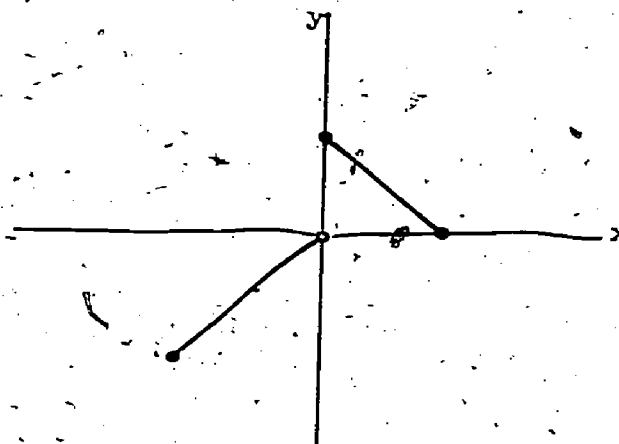
(c)



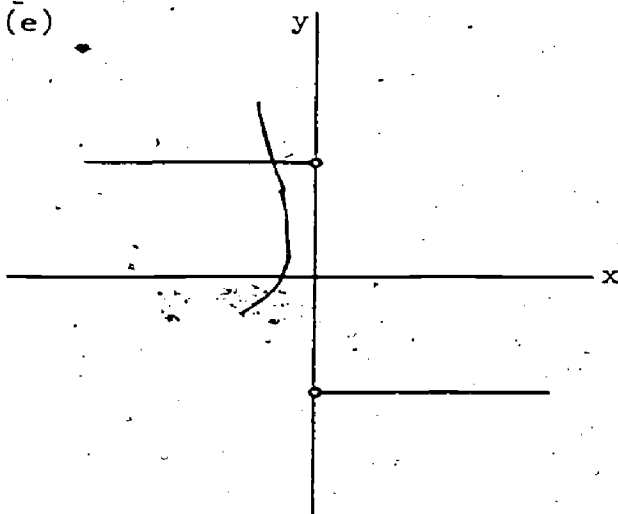
(b)



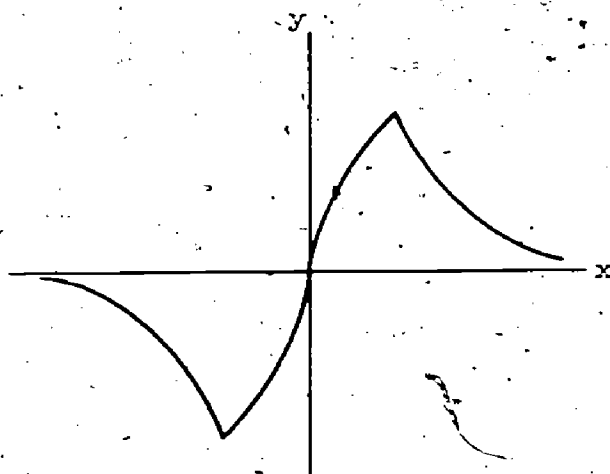
(d)



(e)



(f)



14. Suppose that $f : x \rightarrow f(x)$ is the function whose graph is shown.

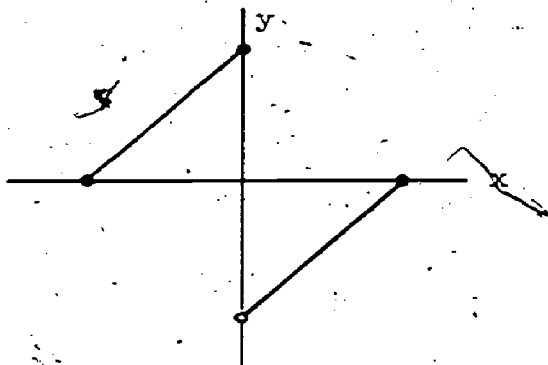
Sketch the graphs of

(a) $g : x \rightarrow 2f(x)$

(b) $g : x \rightarrow f(-x)$

(c) $g : x \rightarrow |f(x)|$

(d) $g : x \rightarrow f(|x|)$



15. A function f is defined by

$$f(x) = \begin{cases} \frac{x}{|x|} & \text{for } x \neq 0, \\ 0 & \text{for } x = 0. \end{cases}$$

Identify this function and sketch its graph.

16. Sketch the graph of each function, specifying its domain and range.

(a) $f : x \rightarrow \sqrt{x^2}$

(g) $f : x \rightarrow \operatorname{sgn} x$

(b) $f : x \rightarrow -|x|$

(h) $f : x \rightarrow \lfloor x \rfloor$

(c) $f : x \rightarrow |1 - x|$

(i) $f : x \rightarrow \frac{\lfloor x \rfloor}{x}$

(d) $f : x \rightarrow 1 - |x|$

(j) $f : x \rightarrow x \lfloor x \rfloor$

(e) $f : x \rightarrow |x| - x$

(k) $f : x \rightarrow |1 - x^2|$

(f) $f : x \rightarrow |x| + |x - 1|$

(l) $f : x \rightarrow |x^2 - 2x - 3|$

(Hint: Consider separately the three possibilities: $x < 0$, $0 \leq x \leq 1$, and $x > 1$.)

Sketch the graphs of the functions in Exercises 17 to 19. For those functions which are periodic, indicate their periods. Indicate those functions which are even or odd.

17. (a) $f : x \rightarrow x - \left[x - \frac{1}{2} \right]$
 (b) $f : x \rightarrow 2x^2 - [2x^2]$
 (c) $f : x \rightarrow 2x^2 - 2[x^2]$
 (d) $f : x \rightarrow 2x^2 - 2[x]^2$
18. (a) $f : x \rightarrow ax - [ax]$, $a > 0$
 (b) $f : x \rightarrow 5x - [2x] - [3x]$
 (c) $f : x \rightarrow x(\sqrt{2} + 1) - [x\sqrt{2}] - [x]$
19. (a) $f : x \rightarrow \frac{1 + \operatorname{sgn} x}{2}$. This function is also called the Heaviside unit function and is designated by $f : x \rightarrow H(x)$.
 (b) $f : x \rightarrow H(x) + H(x - 2)$
 (c) $f : x \rightarrow H(x) \cdot H(x - 2)$
 (d) $f : x \rightarrow (x - 2)^2 \cdot H(x)$
 (e) $f : x \rightarrow H(x) + H(x - 2) + H(x - 4)$
 (f) $f : x \rightarrow H(x^2 - 2)$
 (g) $f : x \rightarrow (\operatorname{sgn} x)(x - 1)^2 + [\operatorname{sgn}(x - 1)]x^2$
20. If f and g are periodic functions of periods m and n , respectively (m, n integers), show that $f + g$ and $f \cdot g$ are also periodic. Give examples to show that the period of $f + g$ can either be greater or less than both of m and n . Repeat the same for the product $f \cdot g$.
21. (a) Can a function be both even and odd?
 (b) What can you say about the evenness or oddness of the product of:
 (1) an even function by an even function?
 (2) an even function by an odd function?
 (3) an odd function by an odd function?
 (c) Show that every function whose domain contains $-x$ whenever it contains x can be expressed as the sum of an even function plus an odd function.

22. Find functions $f(x)$ satisfying

$$f(x) \cdot f(-x) = 1 \quad (\text{called a functional equation}).$$

Suggestion: Use 21(c).

23. Prove that no periodic function other than a constant can be a rational function. (Note: A rational function is the ratio of two polynomial functions.)

A2-2. Composite Functions.

Given two functions f and g with domains whose intersection is non-empty, we can construct new functions by using any of the elementary rational operations--addition, subtraction, multiplication, division--on the given functions. Thus, the sum of f and g is defined to be the function

$$f + g : x \longrightarrow f(x) + g(x)$$

which has for domain those elements contained in the intersection of the domains of f and g . Similarly there are definitions for the difference, product, and quotient of two functions; there is, in fact, a whole algebra of functions, just as there is the familiar algebra of real numbers.

In this algebra of functions there is one operation that has no counterpart in the algebra of numbers: the operation of composition. This operation is best explained by examples.

Let

$$g : x \longrightarrow 2x + 1$$

and

$$f : x \longrightarrow x^2.$$

We observe that

$$g(1) = 3 \quad \text{and} \quad f(3) = 9,$$

$$g(2) = 5 \quad \text{and} \quad f(5) = 25,$$

and, in general, the value of f at $g(x)$ is

$$f(g(x)) = f(2x + 1) = (2x + 1)^2.$$

We have constructed a new function which maps x onto the square of $(2x + 1)$. This function, defined by the mapping $x \longrightarrow f(g(x))$ and denoted by fg , is called a composite of f and g . Hereafter we shall usually represent the value of the function fg by $fg(x)$ rather than $f(g(x))$. Either symbol means the value of f at $g(x)$.*

* The symbol fg denoting the composite of the functions f and g must not be confused with the product of the functions. In this text we distinguish the latter by use of the dot for multiplication; i.e., $f \cdot g$.

An immediate question arises as to the order in which two functions are composed: is the composition of functions a commutative operation; i.e., in general, are $gf(x)$ and $fg(x)$ equal? In the example above we have seen that $fg(2) = f(5) = 25$, and we calculate $gf(2)$:

$$gf(2) = g(4) = 9 \neq fg(2).$$

This counterexample is sufficient to prove that in general $gf(x) \neq fg(x)$. The operation of composition applied to two functions f and g generally produces two different composite functions fg and gf , depending upon the order in which they are composed.

A word of caution must be injected at this point. The number $fg(x)$ is defined only if x is in the domain of g and $g(x)$ is in the domain of f . For example, if

$$f(x) = \sqrt{x} \text{ and } g(x) = 3x - 9,$$

then

$$fg(x) = f(3x - 9) = \sqrt{3x - 9},$$

and the domain of fg is the set of real numbers x for which $3x - 9$ is nonnegative; hence the domain is the set of all $x \geq 3$.

For the other composition of the same functions f and g , we have

$$gf(x) = g(\sqrt{x}) = 3\sqrt{x} - 9$$

which is defined for all nonnegative real numbers x .

We define composition of functions formally.

DEFINITION A2-2. The composite fg of two functions f and g is the function

$$fg : x \longrightarrow fg(x) = f(g(x)).$$

The domain of fg is the set of all elements x in the domain of g for which $g(x)$ is in the domain of f . The operation of forming a composite of two functions is called composition.

The definition may be extended to the composition of three or more functions. Thus, if f , g , and h are functions, one composite is

$$fgh : x \longrightarrow fgh(x) = f(g(h(x)))$$

In order to evaluate $fgh(x)$, we first find $h(x)$, then the value of g at $h(x)$, and finally the value of f at $gh(x)$.

Exercises A2-2

1. Given that $f : x \rightarrow x - 2$ and $g : x \rightarrow x^2 + 1$ for all real x , find

- (a) $f(2) + g(2)$. (e) $f(x) + g(x)$.
 (b) $f(2) \cdot g(2)$. (f) $f(x) \cdot g(x)$.
 (c) $fg(2)$. (g) $fg(x)$.
 (d) $gf(2)$. (h) $gf(x)$.

2. If $f(x) = 3x + 2$ and $g(x) = 5$, find

- (a) $fg(x)$.
 (b) $gf(x)$.

3. If $f(x) = 2x + 1$ and $g(x) = x^2$, find

- (a) $fg(x)$ and $gf(x)$.
 (b) For what values of x , if any, are $fg(x)$ and $gf(x)$ equal?

4. For each pair of functions f and g ; find the composite functions fg and gf and specify the domain (and range, if possible) of each.

(a) $f : x \rightarrow \frac{1}{x}$, $g : x \rightarrow 2x - 6$

(b) $f : x \rightarrow \frac{1}{x}$, $g : x \rightarrow x^2 - 4$

(c) $f : x \rightarrow \frac{1}{x}$, $g : x \rightarrow \sqrt{x}$

(d) $f : x \rightarrow x^2$, $g : x \rightarrow \sqrt{x}$

(e) $f : x \rightarrow x^2$, $g : x \rightarrow \sqrt{4 - x}$

(f) $f : x \rightarrow -x^2 - 1$, $g : x \rightarrow \sqrt{x}$

5. Given that $f(x) = x^2 + 3$ and $g(x) = \sqrt{x + 2}$, solve the equation

$$fg(x) = gf(x)$$

6. Solve problem 5 taking $g(x) = \sqrt{x - 2}$.

7. Describe functions f and g such that gf will equal

(a) $3(x + 2) - 4$.

(d) $\sqrt{x^2 - 4}$.

(b) $(2x - 5)^3$.

(e) $(x^4)^2$.

(c) $\frac{3}{2x - 5}$.

8. For each pair of functions f and g find the composite functions fg and gf and specify the domain (and range, if possible) of each. Also, sketch the graph of each, and give the period (fundamental) of those which are periodic.

(a) $f : x \longrightarrow |x|$, $g : x \longrightarrow \operatorname{sgn}(x - 2)$

(b) $f : x \longrightarrow |x|$, $g : x \longrightarrow 2 \operatorname{sgn}(x - 2) - 1$

9. What can you say about the evenness or oddness of the composite of

(a) an even function of an even function?

(b) an even function of an odd function?

(c) an odd function of an odd function?

(d) an odd function of an even function?

10. If the function f is periodic, what can you say about the periodic character of the composite functions fg and gf assuming these exist and g is an arbitrary function (not periodic)? Illustrate by examples.

11. If the functions f and g are each periodic, then the composite functions fg and gf (assumed to exist) are also periodic. Can the period of either one be less than that of both f and g ?

12. A sequence $a_0, a_1, a_2, \dots, a_n, \dots$, is defined by the equation

$$a_{n+1} = f(a_n), \quad n = 0, 1, 2, 3, \dots,$$

where f is a given function and a_0 is a given number. If $a_0 = 0$ and $f : x \longrightarrow \sqrt{2+x}$; then

$$a_1 = f(a_0) = \sqrt{2}$$

$$a_2 = f(a_1) = ff(a_0) = \sqrt{2 + \sqrt{2}}$$

$$a_3 = f(a_2) = fff(a_0) = \sqrt{2 + \sqrt{2 + \sqrt{2}}}$$

Show that for any n ,

(a) $a_n < 2$.

(b) $a_n > 2 - \frac{1}{2^{n-1}}$, $n > 0$.

A13. If $a_{n+1} = f(a_n)$, $n = 0, 1, 2, \dots$, $a_0 = \mu$, find a_n as a function of μ and n , for the following functions f :

(a) $f : x \longrightarrow a + bx$.

(b) $f : x \longrightarrow x^m$.

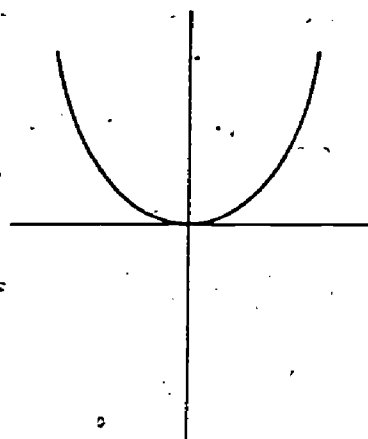
(c) $f : x \longrightarrow \sqrt{|x|}$.

(d) $f : x \longrightarrow \sqrt{1 - x^2}$.

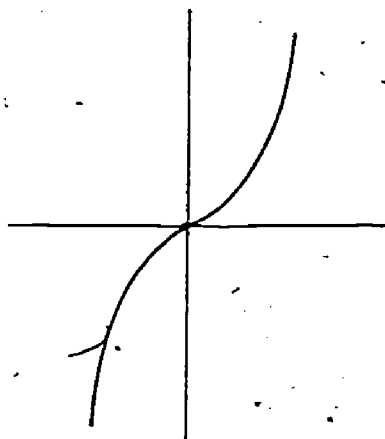
(e) $f : x \longrightarrow (1 - x)^{-1}$.

A2-3. Inverse Functions.

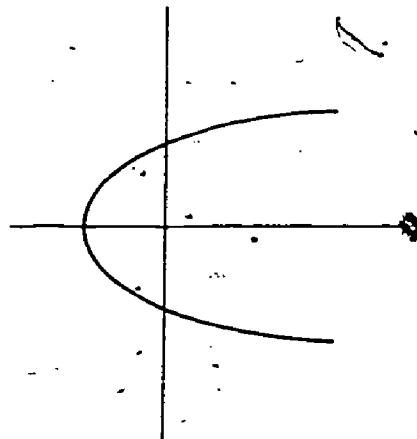
Recall the vertical line test for the graph of a function (Section A2-1) : if every line which is parallel to the y-axis intersects a graph in at most one point, then the graph is that of a function. Thus in Figure A2-3a, (i) and (ii) illustrate graphs of functions, (iii) is the graph of a relation that is not a function.



(i)



(ii)



(iii)

Figure A2-3a

This figure also illustrates an important distinction between two classes of functions: for graph (i) there is at least one line parallel to the x-axis which intersects the graph in more than one point; this is not the case for graph (ii). The latter is typical of a class of functions called one-to-one functions: each element in the domain is mapped into one and only one image in the range, and each element in the range corresponds to one and only one preimage in the domain. In other words, a function of this kind establishes a one-to-one correspondence between the domain and the range of the function.

DEFINITION A2-3a. A function f is one-to-one if whenever $f(x_1) = f(x_2)$, then $x_1 = x_2$:

Note the distinction between Definitions A2-1 and A2-3a. The former states that any function f has the property that if $x_1 = x_2$, then $f(x_1) = f(x_2)$, whereas the latter states that a one-to-one function f is such that $f(x_1) = f(x_2)$ if and only if $x_1 = x_2$.

The class of one-to-one functions is important because for each member of this class we can specify a function that, in a loose way of speaking, undoes the work of the given function. Thus, for example, if f is the

function which maps each real number onto its double, then there is a function g , called the inverse of f , which reverses this mapping and takes each real number onto its half: $f : x \longrightarrow 2x$; $g : y \longrightarrow \frac{1}{2}y$.

DEFINITION A2-3b. If a function $f : x \longrightarrow f(x)$ is one-to-one, then the function $g : f(x) \longrightarrow x$, whose domain is the range of f , is called the inverse of f .

The functions f and g represent the same association but considered from opposite directions; the domain of g is the range of f and the range of g is the domain of f . Furthermore, g is itself one-to-one and its inverse is f .

It is instructive to look at the composites of two functions f and g inverse to one another. If f maps x into y , then g maps y back into x ; in other words, if $y = f(x)$, then $x = g(y)$. Hence,

$$gf(x) = g(y) = x, \text{ for all } x \text{ in the domain of } f,$$

and

$$fg(y) = f(x) = y, \text{ for all } y \text{ in the range of } f.$$

Observe that the restriction of the domain of g to coincide with the range of f is part of the definition of the inverse.

Example A2-3a. Consider the one-to-one function $f : x \longrightarrow 2x - 3$; what is its inverse? Here f is described by the instruction, "Take a number, double it, and then subtract 3." In order to reverse this procedure, we must add 3 and then divide by 2. This suggests that the inverse of f is the function $g : x \longrightarrow \frac{x+3}{2}$. To prove this fact, we must show that g satisfies Definition A2-3b; i.e., show that g maps $f(x)$ into x for all x in the domain of f . By substitution,

$$gf(x) = g(2x - 3) = \frac{(2x - 3) + 3}{2} = x;$$

g is the inverse of f . Furthermore, in the opposite direction,

$$fg(x) = f\left(\frac{x+3}{2}\right) = 2\left(\frac{x+3}{2}\right) - 3 = x$$

for all x in the domain of g . Hence, f is the inverse of the function g , as expected.

The graph of the inverse g of a function f is easily found from the graph of f . If f maps a into b , then g maps b into a . It follows that the point (a,b) is on the graph of f , if and only if (b,a) is on the graph of g . Figure A2-3b shows three points $(1,-3)$, $(2,1)$, and $(4,2)$ on the graph of a function f , and their corresponding points, obtained by interchange of coordinates, on the graph of g :

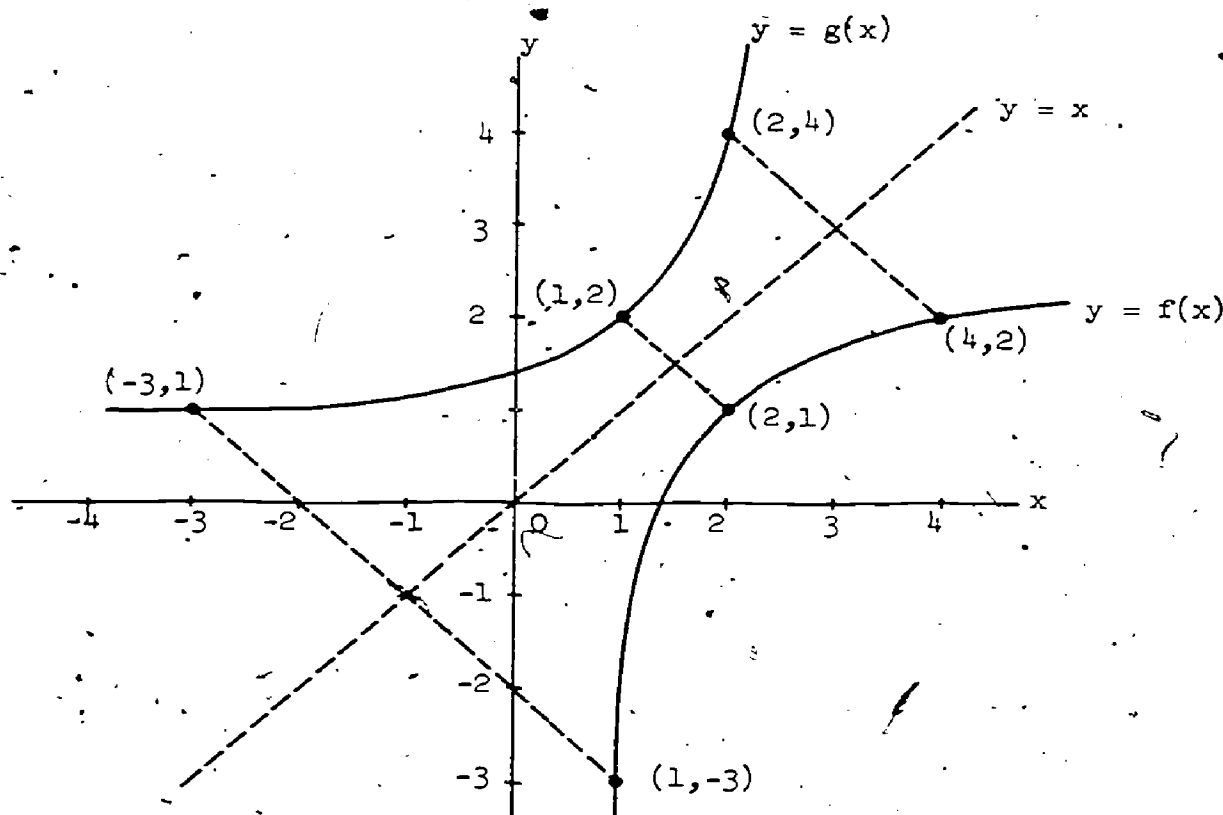


Figure A2-3b

From this figure we see that the points (a,b) and (b,a) are symmetric with respect to the line $y = x$; that is, the line segment determined by these two points is perpendicular to, and bisected by, the line $y = x$. We call (b,a) the reflection of (a,b) in the line $y = x$.

Example A2-3b. Consider the functions $f : x \rightarrow \sqrt{x+2}$, $x \geq -2$, and $g : x \rightarrow x^2 - 2$. The function f is one-to-one; g is not and, hence, cannot be the inverse of f as it stands. By Definition A2-3b, the domain of g must be the range of f , namely, the set of nonnegative real numbers. Hence, the inverse of f is $g : x \rightarrow x^2 - 2$, $x \geq 0$ (Figure A2-3c). The composite functions verify that f and g are inverse to one another:

$$fg : x \rightarrow fg(x) = \sqrt{(x^2 - 2) + 2} = x, \quad x \geq 0;$$

$$gf : x \rightarrow gf(x) = (\sqrt{x + 2})^2 - 2 = x, \quad x \geq -2.$$

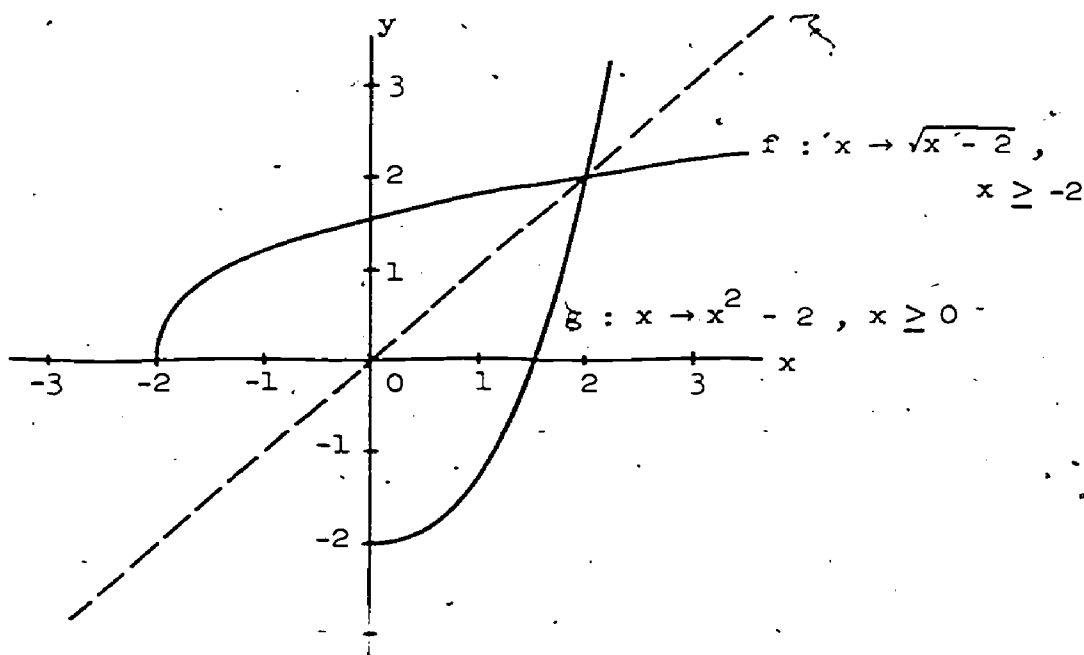


Figure A2-3c

The relationship between the coordinates of a point (a, b) and the coordinates of its reflection (b, a) in the line $y = x$ suggests a formal method for obtaining an equation of the inverse of a given function assuming that the inverse exists.

Example A2-3c. Consider the function

$$f : x \rightarrow y = 3x + 5 \quad \text{for all real } x.$$

If we interchange x and y in the equation

$$(1) \quad y = 3x + 5,$$

we obtain

$$(2) \quad x = 3y + 5.$$

For every pair of numbers (a, b) in the solution set of (1), a pair (b, a) is in the solution set of (2). Hence, (2) is an equation defining implicitly the inverse of the given function f . In order to obtain the explicit form, we solve (2) for y in terms of x and obtain

$$y = \frac{x-5}{3}.$$

The inverse of f is, therefore,

$$g : x \longrightarrow g(x) = \frac{x-5}{3} \quad \text{for all real } x.$$

You should verify the fact that $gf(a) = a$ for any a in the domain of f , and that $fg(b) = b$ for any b in the domain of g (range of f).

Example A2-3d. If the given equation defines a quadratic function, the problem of finding an inverse is more complicated. In the first place, the given function must be restricted to a domain which gives a one-to-one function; in the second place, the technical details of interchanging the variables x and y in the given equation and then solving for y are more involved.

Consider the function

$$f : x \longrightarrow x^2 + 2x + 3$$

whose graph is a parabola with vertex at $(-1, 2)$ and opening upward. If, for example, we restrict f to the domain $\{x : x \geq -1\}$, then we have a function f_1 which is one-to-one and hence has an inverse g_1 . The range of f_1 is $\{y : y = f_1(x) \geq 2\}$, and this will be the domain of g_1 .

We proceed to find a formula defining g_1 . We are given

$$y = x^2 + 2x + 3,$$

and we interchange the variables to obtain

$$x = y^2 + 2y + 3.$$

We now solve for y in the quadratic equation

$$y^2 + 2y + (3 - x) = 0;$$

obtaining

$$y = -1 + \sqrt{x-2} \quad \text{or} \quad y = -1 - \sqrt{x-2}.$$

Which of these formulas defines the function g_1 ? Since y here represents any element in the range of the inverse function, and since the range must be the same set of numbers as the domain of f_1 , we see that $y \geq -1$ is required. Hence

$$y = -1 + \sqrt{x - 2}$$

defines the inverse function

$$g_1 : x \mapsto -1 + \sqrt{x - 2}$$

whose domain is $\{x : x \geq 2\}$. (Note, again, that this is the range of f_1 .)

It is helpful to sketch the graphs of the two inverse functions in order to see more clearly the relationships between their domains and ranges. (See Figure A2-3d.) In fact, if you graph the original function f , you may see more clearly how its domain may be restricted in infinitely many ways to give as many different one-to-one functions, each of which has a unique inverse function.

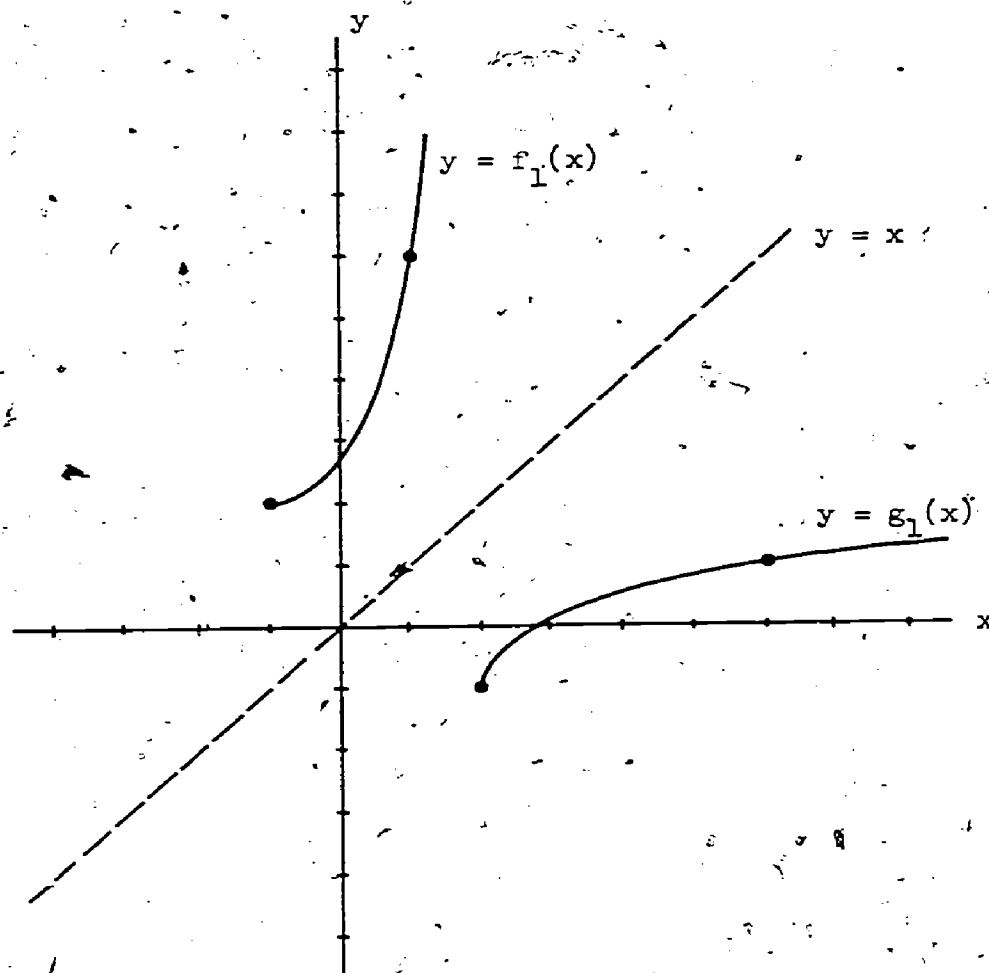
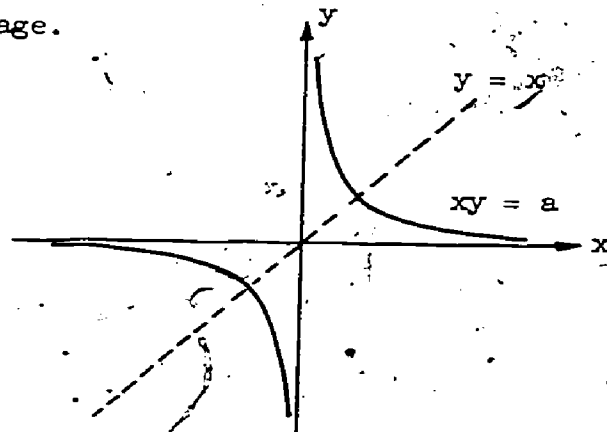


Figure A2-3d

Exercises A2-3

1. What is the reflection of the line $y = f(x) = 3x$ in the line $y = x$? Write an equation defining the inverse of f .
2. Which points are their own reflections in the line $y = x$? What is the graph of all such points?
3. (a) Find the slope of the segment from (a, b) to (b, a) , and prove that the segment is perpendicular to the line $y = x$.
(b) Prove that the segment from (a, b) to (b, a) is bisected by the line $y = x$.
4. What is the reflection of $(1, 1)$ in the line
 - (a) $x = 0$?
 - (b) $y = 0$?
 - (c) $y = -x$?
 - (d) $y = 2$?
 - (e) $x = -3$?
5. Describe any function or functions you can think of which are their own inverses.
6. An equation or an expression (phrase) is said to be symmetric in x and y if the equations or the expressions remain unaltered by interchanging x and y ; e.g., $x^2 + y^2 = 0$, $x^3 + y^3 - 3xy$, $|x - y| = |x + y|$, $x - xy + y$. It follows that graphs of symmetric equations are symmetric about the $y = x$ line. Geometrically, we can consider the line $y = x$ behaving as a mirror, i.e., for any portion of the graph there must also be a portion which is the mirror image.
The equation $x^4 + y^4 = a^4$ is obviously symmetric with respect to the line $y = x$. What other axes of symmetry (mirror type) does it have?



7. The expression

$$a + b + |a - b| + 2c + |a + b + |a - b| - 2c|$$

is obviously symmetric in a and b . Show that it is also symmetric in a and c .

8. Find the inverse of each function.

(a) $f : x \longrightarrow 3x + 6$

(b) $f : x \longrightarrow x^3 - 5$

(c) $f : x \longrightarrow \frac{2}{x} - 3$

9. Which of the following functions have inverses? Describe each inverse by means of a graph or equation and give its domain and range.

(a) $f : x \longrightarrow x^2$

(d) $f : x \longrightarrow [x]$

(b) $f : x \longrightarrow \sqrt{x}$

(e) $f : x \longrightarrow x|x|$

(c) $f : x \longrightarrow |x|$

(f) $f : x \longrightarrow \operatorname{sgn} x$

10. As we have seen, $f : x \longrightarrow x^2$ for all real x does not have an inverse. Do the following:

(a) Sketch graphs of $f_1 : x \longrightarrow x^2$ for $x \geq 0$ and $f_2 : x \longrightarrow x^2$ for $x \leq 0$, and determine the inverses of f_1 and f_2 .

(b) What relationship exists among the domains of f , f_1 , and f_2 ? (f_1 is called the restriction of f to the domain $\{x : x \geq 0\}$, and f_2 is similarly the restriction of f to the domain $\{x : x \leq 0\}$.)

11. (a) Sketch a graph of $f : x \longrightarrow \sqrt{4 - x^2}$ and show that f does not have an inverse.

(b) Divide the domain of f into two parts such that the restriction of f to either part has an inverse.

(c) Write an equation defining each inverse of part (b) and sketch the graphs.

12. Do Problem 11 for $f : x \longrightarrow x^2 - 4x$.

13. Given that $f(x) = 3x - 2$ and $g(x) = -2x + k$, find k such that $fg(x) = gf(x)$. For this value of k , are f and g inverse to one another? Give reasons for your answers.

14. Show that $f : x \longrightarrow x^2 - 4x + 5$ for $x \geq 2$ and $g : x \longrightarrow 2 + \sqrt{x - 1}$ for $x \geq 1$ are inverse to one another by showing that $fg(y) = y$ for all y in the domain of g , and that $gf(x) = x$ for all x in the domain of f .

15. If $f(x) = (2x^3 + 1)^7$, find at least two functions g such that $fg(x) = gf(x)$.

A2-4. Monotone Functions.

If we examine the behavior, for x increasing, of the functions $f : x \rightarrow \sqrt{x}$ and $g : x \rightarrow \sin x$, we note that the values of f increase as x increases, while the values of g are sometimes increasing and sometimes decreasing. Geometrically this means that the graph of f is continually rising as we survey it from left to right (the direction of increasing x), whereas the graph of g , like a wave, is now rising, now falling. The graph of a function may also contain horizontal portions (parallel to the x -axis), where the values of the function remain constant on an interval. A function such as $x \rightarrow [x]$ illustrates this, and also points up the fact that the graph of such a function need not be continuous.

Example A2-4a. The function h , defined by

$$h(x) = \begin{cases} -x^2 & , \quad 0 \leq x \leq 1 ; \\ -1 & , \quad 1 \leq x \leq 2 ; \\ -\frac{x^3}{8} & , \quad 2 \leq x \end{cases}$$

has the graph shown in Figure A2-4a.

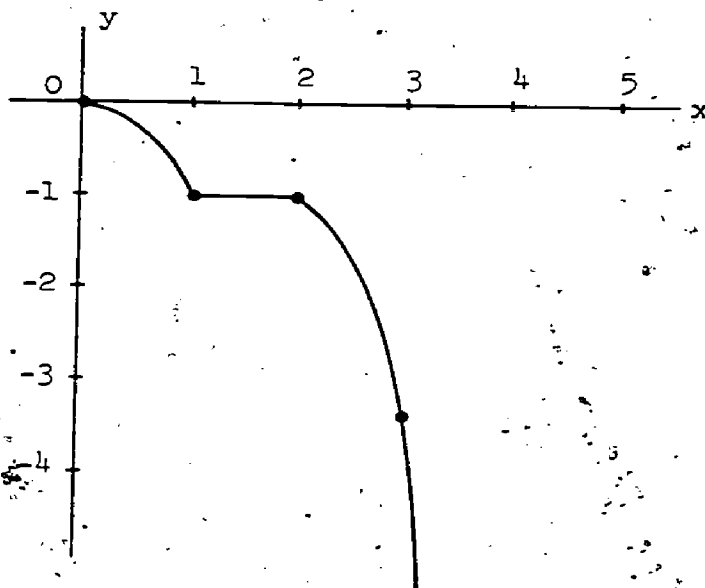


Figure A2-4a

It is easy to see that the function decreases as x increases, except on the interval $[1, 2]$ on which it remains constant. On its entire domain, the values of h are said to be weakly decreasing.

Taken as a class, the increasing, decreasing, weakly increasing, and weakly decreasing functions are called monotone (compare with monotonous) because the changes in the values of the functions as x increases are always in one direction.

DEFINITION A2-4a. Let f be a function defined on an interval I and let $y_1 = f(x_1)$, $y_2 = f(x_2)$ for x_1, x_2 in I . If, for each pair of numbers x_1 and x_2 in I , with $x_1 < x_2$, the corresponding values of y satisfy the inequality

- (1) $y_1 < y_2$, then f is an increasing function;
- (2) $y_1 > y_2$, then f is a decreasing function;
- (3) $y_1 \leq y_2$, then f is a weakly increasing function;
- (4) $y_1 \geq y_2$, then f is a weakly decreasing function.*

Briefly, this definition states that a function which preserves order relations is increasing; a function which reverses order relations is decreasing. Note particularly that an increasing function is a special case of a weakly increasing function; similarly, a decreasing function is a special case of a weakly decreasing function.

DEFINITION A2-4b. A function which is either weakly increasing or weakly decreasing is called monotone. A function which is either increasing or decreasing is called strongly monotone.

For example, the function h of Example A2-4a is monotone over its entire domain and strongly monotone on the closed interval $0 \leq x \leq 1$ as well as on the interval $x \geq 2$.

The graph of a strongly monotone function suggests that the function must be one-to-one, hence must have an inverse.

* In some texts the term "nondecreasing" is used instead of "weakly increasing"; "nonincreasing" is used instead of "weakly decreasing."

THEOREM A2-4. If a function is strongly monotone, then it has an inverse which is strongly monotone in the same sense.

Proof. We treat the case for f increasing; the proof for f decreasing is entirely similar.

By Definitions A2-4a and A2-4b, $f(x_1) = f(x_2)$ if and only if $x_1 = x_2$. Hence, by Definition A2-3a, f is one-to-one, and by Definition A2-3b, f has an inverse

$$g: f(x) \rightarrow x$$

defined for all values $f(x)$ in the range of f .

Finally, g is an increasing function, hence strongly monotone, by Definitions A2-4a and A2-4b.

Example A2-4b. The function

$$f: x \rightarrow x^n,$$

n a natural number, is strongly monotone (increasing) for all real $x \geq 0$. (See Exercise 3-7, number 3.) Hence, f has the inverse function

$$(1) \quad g: x^n \rightarrow x, \quad x \geq 0,$$

which is also an increasing function. For an arbitrary element y in the domain of g , we denote $g(y)$ by $\sqrt[n]{y}$; thus (1) may be rewritten

$$(2) \quad g: y \rightarrow \sqrt[n]{y}, \quad y \geq 0.$$

Comparing (1) and (2), we see that $\sqrt[n]{y}$ is the unique positive solution x of the equation $x^n = y$; we call $\sqrt[n]{y}$ the n -th root of y for all real $y \geq 0$.

If the natural number n is odd, then the function $f: x \rightarrow x^n$ is strongly monotone for all real x , as is its inverse function. This means that every real number has a unique n -th root for n odd. For example, for n odd and a real, $\sqrt[n]{-a^n} = -a$.

If n is even, $f: x \rightarrow x^n$ is decreasing for all real $x \leq 0$, and increasing for all real $x \geq 0$. If f_1 is the restriction of f to the domain $x \geq 0$ and f_2 is the restriction of f to $x \leq 0$, then each of these functions has an inverse, namely

$$g_1 : y \longrightarrow \sqrt[n]{y}$$

and

$$g_2 : y \longrightarrow -\sqrt[n]{y}$$

for n even and all real $y \geq 0$. For n even, the positive n -th root of a nonnegative real number is sometimes called its principal n -th root. The symbol $\sqrt[n]{y}$ always means the principal n -th root.

Exercises A2-4

1. Prove that $f : x \longrightarrow x^2$ for $x \geq 0$ is an increasing function. (Hint: Let $x_1 > x_2 \geq 0$; then $x_1 - x_2 > 0$. From this show that $x_1^2 > x_2^2$.)
2. Which of the following functions are nondecreasing? nonincreasing? decreasing? increasing? In each case the domain is the set of real numbers unless otherwise restricted.

(a) $f_1 : x \longrightarrow c$, c a constant	(h) $f_8 : x \longrightarrow x x $
(b) $f_2 : x \longrightarrow x$	(i) $f_9 : x \longrightarrow x + x $
(c) $f_3 : x \longrightarrow x $	(j) $f_{10} : x \longrightarrow x + x - 1 $
(d) $f_4 : x \longrightarrow [x]$	(k) $f_{11} : x \longrightarrow x - 1 + x - 3 $
(e) $f_5 : x \longrightarrow \operatorname{sgn} x$	(l) $f_{12} : x \longrightarrow f_3 f_4(x)$
(f) $f_6 : x \longrightarrow -x^2$, $x \leq 0$	(m) $f_{13} : x \longrightarrow f_4 f_3(x)$
(g) $f_7 : x \longrightarrow -\sqrt{x}$, $x \geq 0$	
3. For each function in Problem 2 which is not monotone, divide its domain into parts such that the restriction of f to any of these parts gives a monotone or strongly monotone function.

We are given that the function

f_1 is weakly increasing,

f_2 is increasing,

g_1 is weakly decreasing,

g_2 is decreasing,

in a common domain. What is the monotone character, if any, of the following functions:

(a) $f_1 + f_2$

(b) $f_2 + g_1$

(c) $g_1 + g_2$

(d) $g_2 + f_1$

(e) $f_1 \cdot f_2$

(f) $f_2 \cdot g_1$

(g) $g_1 \cdot g_2$

(h) $g_2 \cdot f_1$

(i) $f_1 f_2$

(j) $f_2 f_1$

(k) $f_2 g_1$

(l) $g_1 f_2$

(m) $g_1 g_2$

(n) $g_2 g_1$

(o) $g_2 f_1$

(p) $f_1 g_2$

A2-5. The Circular Functions.

We shall treat the circular (or trigonometric) functions (\sin , \cos , etc.) not so much as the geometrical functions to which you are accustomed, but as purely numerical functions apart from the ideas of geometry. The advantage of approaching the circular functions analytically through the number concept rather than through geometrical concepts is that we can then develop systematic computational techniques for their use and greatly expand their range of application. For example, the property of the circular functions essential for higher analysis is their periodicity, a property to which the considerations of elementary geometry and trigonometry scarcely point. The role of periodicity in our understanding of natural phenomena is profound (see Section 1-3). The circular functions correspond directly to the simplest periodic motions, the turning of a wheel or the motion of a particle transcribing a circle at uniform speed. Yet combinations of these same elementary circular functions can be used to represent the most intricate periodic phenomena. (This is the province of Fourier analysis; the elements of this subject will be within reach when you have completed the calculus.)

The concept of circular function is based upon ideas (like the idea of limit) which are not usually stated precisely before the calculus. Nonetheless, the circular functions are too important to neglect. We shall use them freely, assuming all the properties which are familiar to you from your earlier courses without concern for a logical derivation of these properties. In fact, as the situation warrants, we shall also argue geometrically and intuitively to obtain other properties we may need. Eventually (Section 8-5) we shall be able to define the circular functions purely analytically and derive all the properties used earlier. It is tempting to try to be systematic and to develop the theory of circular functions from the beginning by means of a precise definition, but it is doubtful that such an approach could be made meaningful without the prior intuitively based exploration.

The following fundamental properties* of the circular functions listed below are then to be taken initially as assumptions until we can do better and prove them from a precise definition. It is assumed in the following that the domain of $\sin x$ and $\cos x$ is the set of all real numbers.

*In this text we use radian measure only, as in Properties (1) and (2).

$$(1) \quad \sin x \text{ is periodic with period } 2\pi$$

$$(2) \quad \sin \frac{\pi}{2} = 1$$

$$(3) \quad \sin^2 x + \cos^2 x = 1$$

$$(4) \quad \sin(x + y) = \sin x \cos y + \cos x \sin y$$

$$(5) \quad \sin(-x) = -\sin x$$

This brief list is sufficient for the definition and the derivation of all the other common identities. We list some of the more useful ones below.

$$(6) \quad \cos x = \sin(x + \frac{\pi}{2})$$

$$(7) \quad \cos(-x) = \cos x$$

$$(8) \quad \sin(x \pm y) = \sin x \cos y \pm \cos x \sin y$$

$$(9) \quad \cos(x \pm y) = \cos x \cos y \mp \sin x \sin y$$

$$(10) \quad \sin 2x = 2 \sin x \cos x$$

$$(11) \quad \cos 2x = \cos^2 x - \sin^2 x = 2 \cos^2 x - 1 = 1 - 2 \sin^2 x$$

$$(12) \quad \sin \frac{x}{2} = \pm \sqrt{\frac{1 - \cos x}{2}}$$

$$(13) \quad \cos \frac{x}{2} = \pm \sqrt{\frac{1 + \cos x}{2}}$$

$$(14) \quad \sin x - \sin y = 2 \cos \frac{x+y}{2} \sin \frac{x-y}{2}$$

$$(15) \quad \cos x - \cos y = -2 \sin \frac{x+y}{2} \sin \frac{x-y}{2}$$

$$(16) \quad \tan x = \frac{\sin x}{\cos x} = \frac{1}{\cot x}$$

$$(17) \quad \sec x = \frac{1}{\cos x}$$

$$(18) \quad \csc x = \frac{1}{\sin x}$$

$$(19) \quad 1 + \tan^2 x = \sec^2 x$$

$$(20) \quad \cot^2 x + 1 = \csc^2 x$$

$$(21) \quad \tan(x \pm y) = \frac{\tan x \pm \tan y}{1 \mp \tan x \tan y}$$

$$(22) \quad \tan \frac{x}{2} = \frac{1 - \cos x}{\sin x} = \frac{\sin x}{1 + \cos x}$$

There is one additional fact which may be found useful: a linear combination

$$\alpha \cos x + \beta \sin x$$

can be written in the form

$$A \sin(x + a)$$

where

$$\sin a = \frac{\alpha}{A}, \quad \cos a = \frac{\beta}{A},$$

and

$$A = \sqrt{\alpha^2 + \beta^2}.$$

Exercises A2-5

In the following problems, you may assume properties 1-5 as given on page 304. Also, after any particular identity has been proved it may be used as a true statement in any problem which follows.

1. Prove:

(a) $\sin 0 = 0$.

(d) $\sin \pi = 0$.

(b) $\cos 0 = 1$.

(e) $\cos \pi = -1$.

(c) $\cos \frac{\pi}{2} = 0$.

(f) $\cos(-\frac{\pi}{2}) = 0$.

2. Prove:

(a) $\sin(x + \frac{\pi}{2}) = \cos x$.

(b) $\sin(x + \pi) = -\sin x$.

(c) $\cos(x + \frac{\pi}{2}) = -\sin x$.

(d) $\cos(x + \pi) = -\cos x$.

(e) $\cos x = \cos(-x)$, or equivalently, $\cos(x - \frac{\pi}{2}) = \cos(\frac{\pi}{2} - x)$.

3. Prove:

(a) $\sin(x - y) = \sin x \cos y - \cos x \sin y$.

(b) $\cos(x \pm y) = \cos x \cos y \mp \sin x \sin y$.

(c) $\sin 2t = 2 \sin t \cos t$.

(d) $\cos 2t = 2 \cos^2 t - 1 = 1 - 2 \sin^2 t$.

(e) $\sin 3t = 3 \sin t - 4 \sin^3 t$.

(f) $\cos 3t = 4 \cos^3 t - 3 \cos t$.

(g) $\tan(x \pm y) = \frac{\tan x \pm \tan y}{1 \mp \tan x \tan y}$.

(h) $\sin p - \sin q = 2 \cos \frac{p+q}{2} \sin \frac{p-q}{2}$.

(i) $\cos p - \cos q = -2 \sin \frac{p+q}{2} \sin \frac{p-q}{2}$.

(j) $\tan \frac{t}{2} = \frac{1 - \cos t}{\sin t} = \frac{\sin t}{1 + \cos t}$.

4. Prove:

(a) $\sin \frac{x}{2} = \pm \sqrt{\frac{1 - \cos x}{2}}$.

(b) $\cos \frac{x}{2} = \pm \sqrt{\frac{1 + \cos x}{2}}$.

(c) Explain the significance of the \pm sign in (a) and (b).

5. Determine the numerical values of the following:

(a) $\sin \frac{3\pi}{2}$

(f) $\tan \frac{\pi}{6}$

(b) $\cos \frac{13\pi}{6}$

(g) $\cos \frac{5\pi}{8}$

(c) $\tan \frac{5\pi}{3}$

(h) $\sin \frac{\pi}{16}$

(d) $\sin(-\frac{\pi}{4})$

(i) $\cos^3 \frac{3\pi}{8}$

(e) $\sin \frac{5\pi}{12}$

(j) $\tan \frac{\pi}{5}$

6. Sketch the graph of each of the following functions. For those functions which are periodic, indicate their periods. Indicate those functions which are even or odd.

(a) $f : x \mapsto \sin x \cos x$

(b) $f : x \mapsto \sin \frac{2\pi}{3} x + \cos \pi x$

(c) $f : x \mapsto \sin x + \cos x \sqrt{2}$

(d) $f : x \mapsto \sin 2x$

(e) $f : x \longrightarrow |\sin 3\pi x|$

(f) $f : x \longrightarrow \cos^2 \pi x$

(g) $f : x \longrightarrow 2 \sin x^2$

7. (a) What is the period of $\sin ax$, $a \neq 0$?(b) What is the period of $\sin(ax + b)$, $a \neq 0$?(c) For what values of a and b is the function odd? even?8. For each pair of functions f and g find the composite functions fg and gf and specify the domain and range (if possible) of each. Also, sketch the graph of each, and give the period (fundamental) of those which are periodic.

(a) $f : x \longrightarrow \sin \pi x$; $g : x \longrightarrow \operatorname{sgn} x$.

(b) $f : x \longrightarrow \sin ax$, $a > 0$; $g : x \longrightarrow \operatorname{sgn} x$.

(c) $f : x \longrightarrow 3|\sin x|$; $g : x \longrightarrow [x]$.

(d) $f : x \longrightarrow \max\{\sin \pi x, x - [x]\}$; $g : x \longrightarrow \min\{x, x^2\}$.

9. Solve for all x in the interval $0 \leq x < 2\pi$:

$$\sin^m x + \cos^m x = 1, \quad (m \text{ an integer } > 2).$$

10. If (3) on page 304 is replaced by $\cos x = \sin(x + \frac{\pi}{2})$, then using (1) to (5), show that $\cos^2 rx + \sin^2 rx = 1$, where r is rational.

11. Show how to solve the cubic equation

$$4x^3 - 3x = a \quad (|a| \leq 1)$$

trigonometrically.

A2-6. Polar Coordinates.

A fundamental link between the algebra of numbers and the geometry of points was forged by Rene Descartes (1596 - 1650) when he introduced the notion of representing a point in the plane by means of an ordered pair of numbers. Beautiful in its simplicity, this concept paved the way for the development of coordinate geometry and calculus. By introducing a pair of perpendicular number lines (coordinate axes), Descartes was able to assign to each point in the plane a unique pair of real numbers. We call these numbers the (rectangular) cartesian coordinates of the point; you have used them ever since you first began to represent equations by their graphs.

Other coordinate systems have been invented since Descartes' time because they are better adapted to treat some problems which are awkward to handle in cartesian coordinates. We consider here the polar coordinate system and a few examples of its use.

We suppose that we already have a rectangular (cartesian) coordinate system in the plane. We locate a point P in the plane by polar coordinates, an ordered pair of real numbers $(r, \theta)^*$ where $|r|$ is the length of the segment \overline{OP} (sometimes called the radius vector) and θ is the direction angle made by \overline{OP} with the positive x-axis (polar axis) (Figure A2-6a).

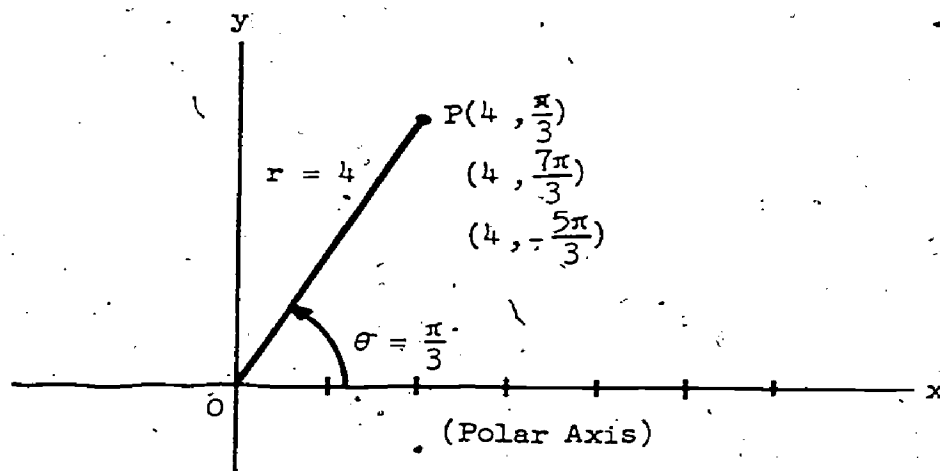


Figure A2-6a

* r is sometimes called the radial coordinate and θ the polar angle or azimuth.

There are infinitely many such angles for each point P ; if θ is one angle, then $\theta \pm 2n\pi$ ($n = 1, 2, 3, \dots$) are the others. Thus, a point may be identified by infinitely many pairs of polar coordinates. For example (Figure A2-6a), point P with polar coordinates $(4, \frac{\pi}{3})$, also has coordinates $(4, \frac{7\pi}{3})$, $(4, -\frac{5\pi}{3})$, and, in general, $(4, \frac{\pi}{3} + 2n\pi)$ for any integer n . The pole (origin) is a special case: to it we assign as polar coordinates any pair $(0, \theta)$, θ any real number.

When we assign polar coordinates to locate a point, it is customary to allow r also to be negative. For $r > 0$, the point $(-r, \theta)$ is located symmetrically to the point (r, θ) with respect to the origin (Figure A2-6b); it has coordinates $(r, \theta + \pi)$ also.

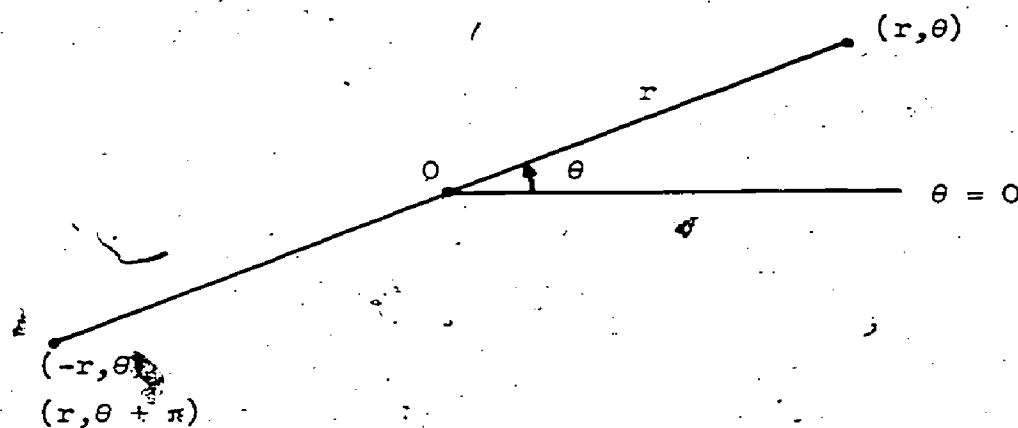


Figure A2-6b

In a cartesian coordinate system every point in the plane has a unique pair of coordinates (x, y) . In a polar coordinate system, by contrast, this is not true; a given point in the plane does not have a unique representation (r, θ) in polar coordinates (see point P in Figure A2-6a). In both coordinate systems, however, a given pair of coordinates specifies a unique point in the plane.

A relation between x and y may be represented by a graph in a cartesian coordinate plane. A relation in r and θ may be represented by a graph in a polar coordinate system; a point lies on the graph if and only if it has at least one coordinate pair which satisfies the given relation. We discuss the graphs of a few functions defined by equations in polar coordinates.

The graph of the equation

$$r = c$$

contains all points (c, θ) , θ any real number; it is a circle of radius $|c|$ having its center at the pole. The equation $r = -c$ describes the same circle.

The points for which

$$\theta = c$$

lie on the line passing through the pole which forms an angle of c radians with the polar axis; each point of the line has coordinates (r, c) for some real r . For r positive, the points form the ray in direction θ , for r negative, the ray has direction $\theta + \pi$. The line has infinitely many equations $\theta = c + n\pi$, n an integer.

The circular functions of θ are especially conveniently represented in polar coordinates because the entire graph is traced out in one period. We shall illustrate a procedure for sketching a graph of such a function using polar coordinate graph paper. Note that the function specifies the graph; a function, however, cannot be recovered from its graph in polar coordinates.

Example A2-6a. Sketch a graph of the function defined by

$$r = 4 \cos \theta$$

Since r is a function of θ , we consider values of θ and calculate the corresponding values of r . We know that the cosine increases from the value 0 at $\theta = -\frac{\pi}{2}$ to 1 at $\theta = 0$ and then decreases to 0 at $\theta = \frac{\pi}{2}$. Hence, in this interval, r increases from 0 to 4 and then decreases to 0. Since $\cos(\theta + \pi) = -\cos \theta$, the point $(4 \cos(\theta + \pi), \theta)$ is the same as $(-4 \cos \theta, \theta)$, and the curve for $-\frac{\pi}{2} \leq \theta \leq \frac{\pi}{2}$ is the entire graph.

To sketch the graph of the function, we calculate r for a few convenient values of θ ($-\frac{\pi}{3}$, $-\frac{\pi}{4}$, $\frac{\pi}{6}$, etc.), locate the corresponding points on polar coordinate paper, and sketch the graph (Figure A2-6c); it appears to be a circle and we shall presently verify that it is.

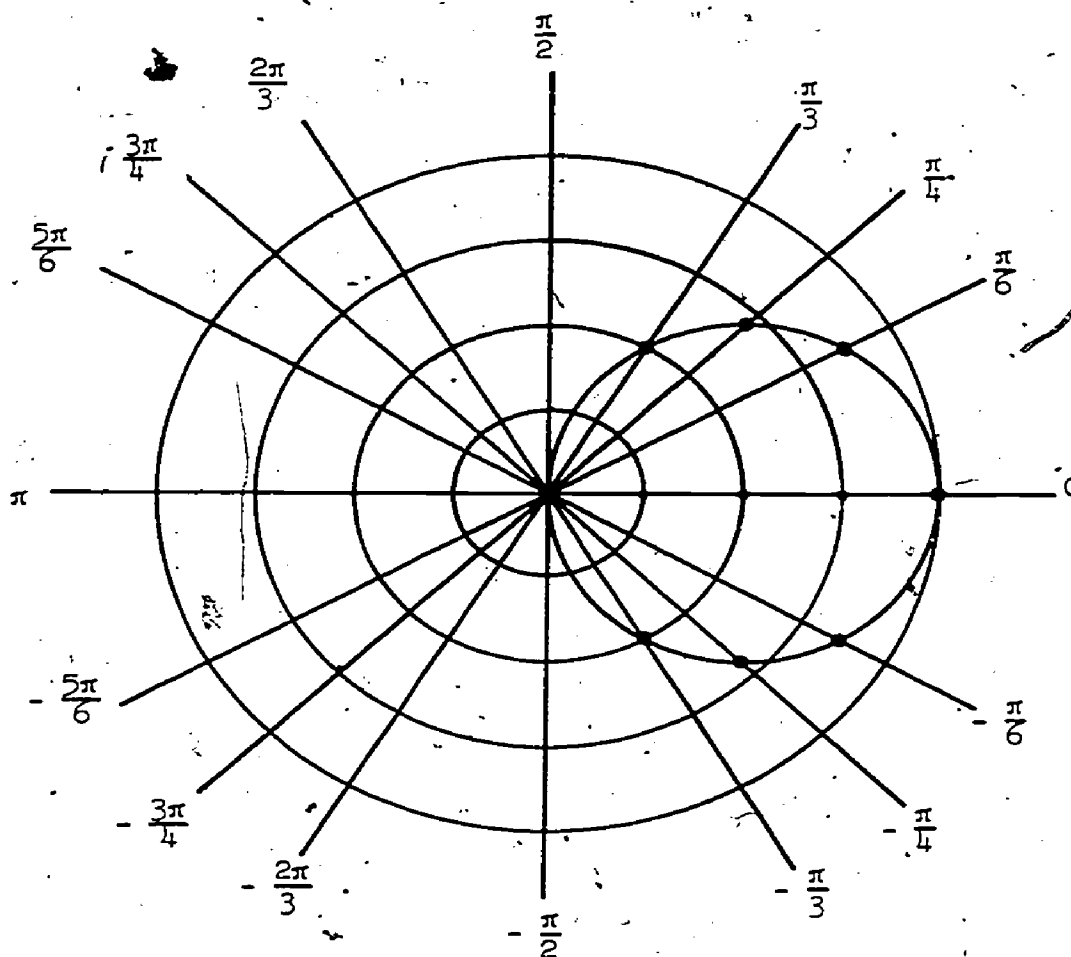


Figure A2-6c

Since each point P in the plane has both rectangular and polar coordinates (Figure A2-6d), for $r > 0$, we have from the trigonometric functions of angles

$$(1) \quad x = r \cos \theta, \quad y = r \sin \theta.$$

We leave it to you to verify that equations (1) hold for $r \leq 0$. Thus the rectangular and polar coordinates of each point in the plane are related

by (1). It follows that

$$(2) \quad x^2 + y^2 = r^2.$$

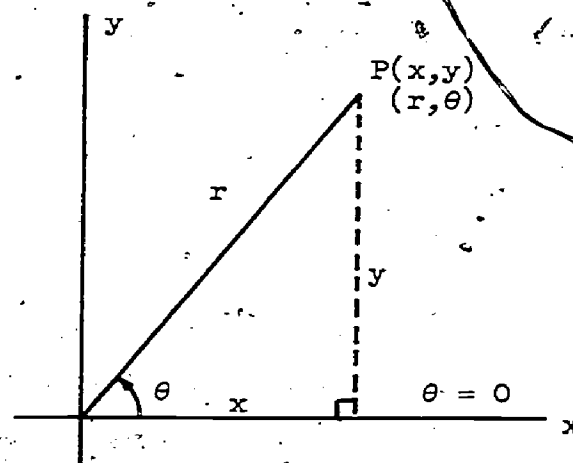


Figure A2-6d

Now we re-examine the function defined by $r = 4 \cos \theta$. (Example A2-4a) and prove that its graph is a circle. We shall do so by transforming the given equation into an equation involving rectangular coordinates x and y . Now the given equation $r = 4 \cos \theta$ has the same graph as the equation

$$(3) \quad r^2 = 4r \cos \theta,$$

for if $r \neq 0$, we may divide both members of the equation by r to obtain the given equation; $r = 0$ corresponds to the fact that the pole is on both graphs. This may not be immediately obvious since only certain pairs of coordinates representing the pole will satisfy the equation $r = 4 \cos \theta$. For example, both $(0,0)$ and $(0, \frac{\pi}{2})$ represent the pole, yet only the latter of these pairs satisfies $r = 4 \cos \theta$.

We use (1) and (2) to obtain from (3) that

$$x^2 + y^2 = 4x$$

or

$$(x - 2)^2 + y^2 = 4,$$

an equivalent equation in rectangular coordinates. We recognize this as an equation of the circle with center at $(2,0)$ and radius 2, verifying the graph in Figure A2-6c.

Example A2-6b. Find an equation in polar coordinates of the curve whose equation in cartesian coordinates is $(x^2 + y^2)^2 = a^2(x^2 - y^2)$.

Applying Equations (1) and (2), we have

$$\begin{aligned} r^4 &= a^2 r^2 (\cos^2 \theta - \sin^2 \theta) \\ &= a^2 r^2 \cos 2\theta. \end{aligned}$$

This is equivalent to

$$r^2 = 0 \text{ (the pole) and } r^2 = a^2 \cos 2\theta.$$

Since $r^2 = a^2 \cos 2\theta$ is satisfied by $(0, \frac{\pi}{4})$, a set of polar coordinates for the pole, we see that $r^2 = 0$ contributes no points not in the graph of $r^2 = a^2 \cos 2\theta$. Hence, the latter is an equation in polar form which is the equivalent of the given one. The graph of this equation is called the

Lemniscate of Bernoulli and is displayed in Figure A5a.*

In Section 12-4 we use an equation of an ellipse in polar form; here we develop an equation which, for suitable choice of a parameter, will represent either a parabola, an ellipse, or a hyperbola. For this purpose we need the definition of these curves (conic sections) in terms of focus, directrix, and eccentricity. Every conic section (other than the circle) may be defined to be the set (locus) of all points P such that the ratio of the distance between P and a fixed point F (the focus) to the distance between P and a fixed line ℓ (the directrix) is a positive constant e , called the eccentricity of the conic section. If $e = 1$ the conic section is a parabola, if $0 < e < 1$ it is an ellipse, and if $e > 1$ it is a hyperbola.

In order to derive an equation in polar coordinates of a conic section, it is convenient to place the focus F at the pole (origin) and the directrix ℓ perpendicular to the extension of the polar axis at distance $p > 0$ from the pole, as shown in Figure A2-6e. (Other orientations are possible; see Exercises A2-6, Nos. 8-10.) Point P is any point of the conic section.

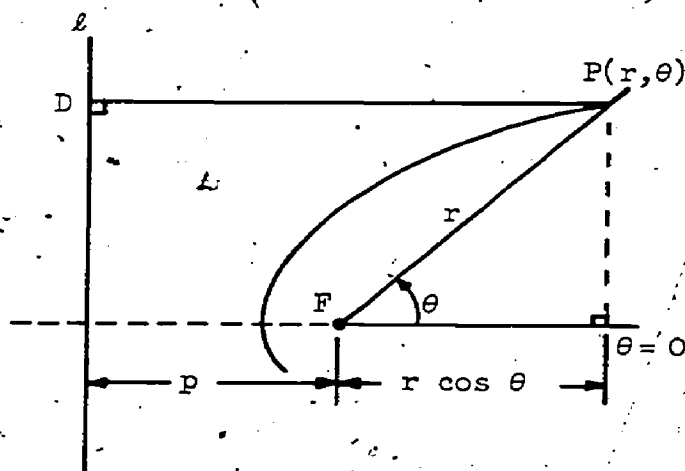


Figure A2-6e

We let (r, θ) be any pair of polar coordinates of P for which $r > 0$; then $FP = r$ and $DP = p + r \cos \theta$ (Figure A2-6e). The definition of the conic sections requires that $\frac{FP}{DP} = e$ or $\frac{r}{p + r \cos \theta} = e$. Solving for r we obtain

$$(4) \quad r = \frac{ep}{1 - e \cos \theta}$$

which we take to be the standard form of the polar equation of conic sections having focus and directrix oriented as in Figure A2-6e. From Equation (4), if $e \leq 1$ (ellipse or parabola), then $r > 0$; if $e > 1$ (hyperbola), r may be negative and these values give us the branch of the hyperbola lying to the left of the directrix.

*This curve is defined as the set (locus) of points P such that the product of the distances of P from two fixed points is the square of half the distance between the two fixed points.

Example A2-6c. Describe and sketch the graph of the equation

$$r = \frac{16}{5 - 3 \cos \theta}$$

We may put this equation in the standard form

$$\begin{aligned} r &= \frac{\frac{16}{5}}{1 - \frac{3}{5} \cos \theta} \\ &= \frac{\frac{3}{5} \cdot \frac{16}{3}}{1 - \frac{3}{5} \cos \theta} \end{aligned}$$

from which $e = \frac{3}{5}$, and $p = \frac{16}{3}$. Since $e < 1$, the graph is an ellipse with focus F_1 at the pole and major axis on the polar axis. By giving θ the values 0 and π , we find the ends of the major axis to be $(8, 0)$ and $(2, \pi)$. Thus the length of the major axis is 10, the center of the ellipse is the point $(3, 0)$, and the other focus is the point $F_2(6, 0)$. Since $p = \frac{16}{3}$ (the distance between a focus and corresponding directrix of the ellipse), the equation of the directrix ℓ_1 corresponding to the focus at the pole is $r \cos \theta = -\frac{16}{3}$ (see Exercises A2-6, No. 6a.), and the equation of the directrix ℓ_2 corresponding to $F_2(6, 0)$ is $r \cos \theta = \frac{34}{3}$. When $\theta = \frac{\pi}{2}$, then $r = \frac{16}{5}$, and we have the point $(\frac{16}{5}, \frac{\pi}{2})$ at one end of the focal chord (latus rectum) through F_1 . The other endpoint has polar coordinates $(\frac{16}{5}, \frac{3\pi}{2})$; these points help us to sketch the ellipse as shown in Figure A2-6f.

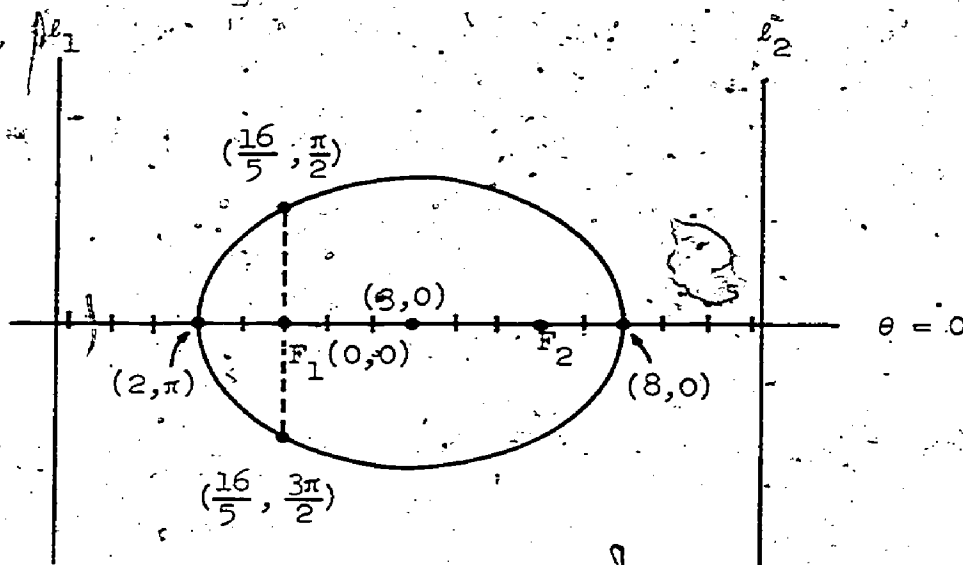


Figure A2-6f

Exercises A2-6

1. Find all polar coordinates of each of the following points:

(a) $(6, \frac{\pi}{4})$	(c) $(6, -\frac{\pi}{4})$
(b) $(-6, \frac{\pi}{4})$	(d) $(-6, -\frac{\pi}{4})$
2. Find rectangular coordinates of the points in Exercise 1.
3. Find polar coordinates of each of the following points given in rectangular coordinates:

(a) $(4, -4)$	(e) $(-3, 0)$
(b) $(\frac{-3\sqrt{3}}{2}, \frac{3}{2})$	(f) $(-3, 4)$
(c) $(-2, -2\sqrt{3})$	(g) $(-\sqrt{3}, 1)$
(d) $(0, -10)$	(h) $(\sqrt{2}, -\sqrt{2})$
4. Given the cartesian coordinates (x, y) of a point, formulate unique polar coordinates (r, θ) for $0 \leq \theta \leq \pi$. (Hint: use $\arccos \frac{x}{r}$.)
5. Determine the polar coordinates of the three vertices of an equilateral triangle if a side of the triangle has length L , the centroid of the triangle coincides with the pole, and one angular coordinate of a vertex is θ_1 radians.
6. Find equations in polar coordinates of the following curves:

(a) $x = c$, c a constant.
(b) $y = c$, c a constant.
(c) $ax + by = c$.
(d) $x^2 + (y - k)^2 = k^2$.
(e) $y^2 = 4ax$.
(f) $x^2 - y^2 = a^2$.
7. Find equations in rectangular coordinates of the following curves:

(a) $r = a$.
(b) $r \sin \theta = -5$.
(c) $r = 2a \sin \theta$.
(d) $r = \frac{1}{1 - \cos \theta}$.
(e) $r = 2 \tan \theta$.

8. Derive an equation in polar coordinates for conic sections with a focus at the pole and directrix perpendicular to the polar axis and p units to the right of the pole.
9. Repeat Number 8 if the directrix is parallel to the polar axis and p units above the focus at the pole.
10. Repeat Number 8 if the directrix is parallel to the polar axis and p units below the focus at the pole.
11. Discuss and sketch each of the following curves in polar coordinates. (See Example A2-6c and Numbers 8, 9, 10.)

(a) $r = \frac{8}{1 - \cos \theta}$

(b) $r = \frac{12}{1 - 3 \cos \theta}$

(c) $r = \frac{36}{5 - 4 \sin \theta}$

(d) $r = \frac{16}{5 + 3 \sin \theta}$

(e) $r \sin \theta = 1 - r$

12. Certain types of symmetry of curves in polar coordinates are readily detected. For example, a curve is symmetric about the pole if the equation is unchanged when r is replaced by $-r$. What kind of symmetry occurs if an equation is unchanged when

(a) θ is replaced by $-\theta$?

(b) θ is replaced by $\pi - \theta$?

(c) r and θ are replaced by $-r$ and $-\theta$, respectively?

(d) θ is replaced by $\pi + \theta$?

13. Without actually sketching the graphs, describe the symmetries of the graphs of the following equations:

(a) $r^2 = 4 \sin 2\theta$.

(b) $r(1 - \cos \theta) = 10$.

(c) $r = \cos^2 2\theta$.

14. Sketch the following curves in polar coordinates:

(a) $r = a\theta$. (d) $r = a^2 \sin^2 \theta \cos^2 \theta$.

(b) $r = a(1 - \cos \theta)$. (e) $r\theta = a$.

(c) $r = a \sin 2\theta$.

15. In each of the following, find all points of intersection of the given pairs of equations. (Recall that the polar representation of a point is not unique.)

(a) $r = 2 - 2 \sin \theta$, $r = 2 - 2 \cos \theta$.

(b) $r = -2 \sin 2\theta$, $r = 2 \cos \theta$.

(c) $r = 4(1 + \cos \theta)$, $r(1 - \cos \theta) = 3$.

Appendix 3

MATHEMATICAL INDUCTION

A3-1. The Principle of Mathematical Induction.

The ability to form general hypotheses in the light of a limited number of facts is one of the most important signs of creativeness in a mathematician. Equally important is the ability to prove these guesses. The best way to show how to guess at a general principle from limited observations is to give examples.

Example A3-1a. Consider the sums of consecutive odd integers.

$$1 = 1$$

$$1 + 3 = 4$$

$$1 + 3 + 5 = 9$$

$$1 + 3 + 5 + 7 = 16$$

$$1 + 3 + 5 + 7 + 9 = 25$$

Notice that in each case the sum is the square of the number of terms.

Conjecture: The sum of the first n odd positive integers is n^2 .
(This is true. Can you show it?)

Example A3-1b. Consider the following inequalities:

$$1 < 100, 2 < 100, 3 < 100, 4 < 100, 5 < 100, \text{ etc.}$$

Conjecture: All positive integers are less than 100. (False, of course.)

Example A3-1c. Consider the number of complex zeros, including the repetitions, for polynomials of various degrees.

Zero degree: a_0 , no zeros ($a_0 \neq 0$) .

First degree: $a_1x + a_0$, one zero at $x = \frac{-a_0}{a_1}$.

Second degree: $a_2x^2 + a_1x + a_0$, two zeros at

$$x = \frac{-a_1 \pm \sqrt{a_1^2 - 4a_0a_2}}{2a_2}$$

Conjecture: Every polynomial of degree n has exactly n complex zeros when repetitions are counted. (True.)

Example A3-1d. Observe the operations necessary to compute the roots from the coefficients in Example A3-1c.

Conjecture: The zeros of a polynomial of degree n can be given in terms of the coefficients by a formula which involves only addition, subtraction, multiplication, division, and the extraction of roots. (False.)

Example A3-1e. Take any even number except 2 and try to express it as the sum of as few primes as possible:

$$4 = 2 + 2 , 6 = 3 + 3 , 8 = 3 + 5 , 10 = 5 + 5 ,$$

$$12 = 5 + 7 , 14 = 7 + 7 , \text{ etc.}$$

Conjecture: Every even number but 2 can be expressed as the sum of two primes. (As yet, no one has been able to prove or disprove this conjecture.)

Common to all these examples is the fact that we are trying to assert something about all the members of a sequence of things: the sequence of odd integers, the sequence of positive integers, the sequence of degrees of polynomials, the sequence of even numbers greater than 2. The sequential character of the problems naturally leads to the idea of sequential proof. If we know something is true for the first few members of the sequence, can we use that result to prove its truth for the next member of the sequence? Having done that, can we now carry the proof on to one more member? Can we repeat the process indefinitely?

Let us try the idea of sequential proof on Example A3-1a. Suppose we know that for the first k odd integers $1, 3, 5, \dots, 2k-1$,

$$(1) \quad 1 + 3 + 5 + \dots + (2k-1) = k^2 ,$$

can we prove that upon adding the next higher odd number $(2k + 1)$ we obtain the next higher square? From (1) we have at once by adding $2k + 1$ on both sides,

$$[1 + 3 + 5 + \dots + (2k - 1)] + (2k + 1) = k^2 + (2k + 1) = (k + 1)^2.$$

It is clear that if the conjecture of Example A3-1a is true at any stage then it is true at the next stage. Since it is true for the first stage, it must be true for the second stage, therefore true for the third stage, hence the fourth, the fifth, and so on forever.

Example A3-1f. In many good toy shops there is a puzzle which consists of three pegs and a set of graduated discs as depicted in Figure A3-1a. The problem posed is to transfer the pile of discs from one peg to another under the following rules:

1. Only one disc at a time may be transferred from one peg to another.
2. No disc may ever be placed over a smaller disc.

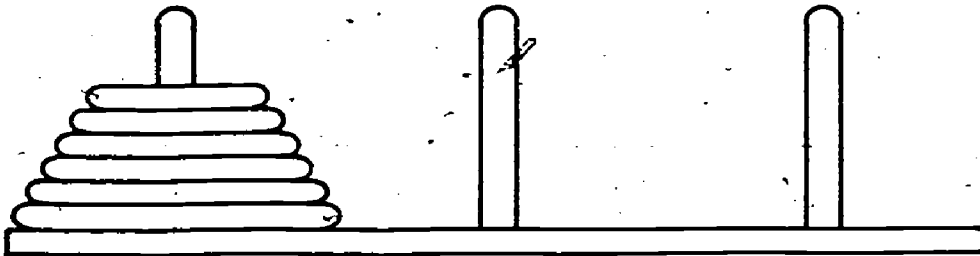


Figure A3-1a

Two questions arise naturally: Is it possible to execute the task under the stated restrictions? If it is possible, how many moves does it take to complete the transfer of the discs? If it were not for the idea of sequential proof, one might have difficulty in attacking these questions.

As it is, we observe that there is no problem in transferring one disc.

If we have to transfer two discs, we transfer one, leaving a peg free for the second disc; we then transfer the second disc and cover with the first.

If we have to transfer three discs, we transfer the top two, as above. This leaves a peg for the third disc to which it is then moved, and the first two discs are then transferred to cover the third disc.

The pattern has now emerged. If we know how to transfer k discs, we can transfer $k + 1$ in the following way. First, we transfer k discs leaving the $(k + 1)$ -th disc free to move to a new peg; we move the $(k + 1)$ -th disc and then transfer the k discs again to cover it. We see then that it is possible to move any number of graduated discs from one peg to another without violating the rules (1) and (2), since knowing how to move one disc, we have a rule which tells us how to transfer two, and then how to transfer three, and so on.

To determine the smallest number of moves it takes to transfer a pile of discs, we observe that no disc can be moved unless all the discs above it have been transferred, leaving a free peg to which to move it. Let us designate by m_k the minimum number of moves needed to transfer k discs. To move the $(k + 1)$ -th disc, we first need m_k moves to transfer the discs above it to another peg. After that we can transfer the $(k + 1)$ -th disc to the free peg. To move the $(k + 2)$ -th disc (or to conclude the game if the $(k + 1)$ -th disc is last) we must now cover the $(k + 1)$ -th disc with the preceding k discs; this transfer of the k discs cannot be accomplished in less than m_k moves. We see then that the minimum number of moves for $k + 1$ discs is

$$m_{k+1} = 2m_k + 1.$$

This is a recursive expression for the minimum number of moves, that is, if the minimum is known for a certain number of discs, we can calculate the minimum for one more disc. In this way, we have defined the minimum number of sequential moves: by adding one disc we increase the necessary number of moves to one more than twice the preceding number. It takes one move to move one disc, therefore it takes three moves to move two discs, and so on.

Let us make a little table (Table I).

Table 1

k	1	2	3	4	5	6	7
m_k	1	3	7	15	31	63	127

k = number of discs

m_k = minimum number of moves

Upon adding a disc we roughly double the number of moves. This leads us to compare the number of moves with the powers of two: 1, 2, 4, 8, 16, 32, 64, 128, ...; and we guess that $m_k = 2^k - 1$. If this is true for some value k , we can easily see that it must be true for the next, for we have

$$\begin{aligned}
 m_{k+1} &= 2m_k + 1 \\
 &= 2(2^k - 1) + 1 \\
 &= 2^{k+1} - 2 + 1 \\
 &= 2^{k+1} - 1,
 \end{aligned}$$

and this is the value of $2^n - 1$, for $n = k + 1$. We know that the formula for m_k is valid when $k = 1$, but now we can prove in sequence that it is true for 2, 3, 4, and so on.

According to persistent rumor, there is a puzzle of this kind in a most holy monastery hidden deep in the Himalayas. The puzzle consists of 64 discs of pure beaten gold and the pegs are diamond needles. The story relates that the game of transferring the discs has been played night and day by the monks since the beginning of the world, and has yet to be concluded. It also has been said that when the 64 discs are completely transferred, the world will come to an end. The physicists say the earth is about four billion years old, give or take a billion or two. Assuming that the monks move one disc every second and play in the minimum number of moves, is there any cause for panic? (Cf. Ball, W. W., Mathematical Recreations. New York: Macmillan Co., 1947; p. 303 ff.)

The principle of sequential proof, stated explicitly, is this (First Principle of Mathematical Induction): Let A_1, A_2, A_3, \dots be a sequence of assertions, and let H be the hypothesis that all of these are true. The hypothesis H will be accepted as proved if

1. There is a general proof to show that if any assertion A_k is true, then the next assertion A_{k+1} is true;

2. There is a special proof to show that A_1 is true.

If there are only a finite number of assertions in the sequence, say ten, then we need only carry out the chain of ten proofs explicitly to have a complete proof. If the assertions continue in sequence endlessly, as in Example 1, then we cannot possibly verify directly every link in the chain of proof. It is just for this reason--in effect that we can handle an infinite chain of proof without specifically examining every link--that the concept of sequential proof becomes so valuable. It is, in fact, at the heart of the logical development of mathematics.

Through an unfortunate association of concepts this method of sequential proof has been named "mathematical induction." Induction, in its common English sense, is the guessing of general propositions from a number of observed facts. This is the way one arrives at assertions to prove. "Mathematical induction" is actually a method of deduction or proof and not a procedure of guessing, although to use it we ordinarily must have some guess to test. This usage has been in the language for a long time, and we would gain nothing by changing it now. Let us keep it then, and remember that mathematical usage is special and often does not resemble in any respect the usage of common English.

In Example A3-1a, above, the assertion A_n is

$$1 + 3 + 5 + \dots + (2n - 1) = n^2.$$

We proved first, that if A_k is true (that is, if the sum of the first k odd numbers is k^2) then A_{k+1} is true, so that the sum of the first $k + 1$ odd numbers is $(k + 1)^2$. Second, we observed that A_1 is true: $1 = 1^2$. These two steps complete the proof.

Mathematical induction is a method of proving a hypothesis about a list or sequence of assertions. Unfortunately it doesn't tell us how to make the hypothesis in the first place. In the example just considered, it was easy to guess from a few specific instances that the sum of the first n odd numbers is n^2 , but the next problem (Example A3-1g) may not be so obvious.

Example A3-1g. Consider the sum of the squares of the first n positive integers,

$$1^2 + 2^2 + 3^2 + \dots + n^2.$$

We find that when $n = 1$, the sum is 1; when $n = 2$, the sum is 5; when

$n = 3$, the sum is 14; and so on. Let us make a table of the first few values (Table II).

Table II

n	1	2	3	4	5	6	7	8
sum	1	5	14	30	55	91	140	204

Though some mathematicians might be immediately able to see a formula that will give us the sum, most of us would have to admit that the situation is obscure... We must look around for some trick to help us discover the pattern which is surely there; what we do will therefore be a personal, individual matter. It is a mistake to think that only one approach is possible.

Sometimes experience is a useful guide. Do we know the solutions to any similar problems? Well, we have here the sum of a sequence, and Example A3a also dealt with the sum of a sequence: the sum of the first n odd numbers is n^2 . Consider the sum of the first n integers themselves (not their squares)--what is

$$1 + 2 + 3 + \dots + n ?$$

This seems to be a related problem, and we can solve it with ease. The terms form an arithmetic progression in which the first term is 1 and the common difference is also 1; the sum, by the usual formula, is therefore

$$\frac{n}{2}(n+1) = \frac{1}{2}n^2 + \frac{1}{2}n.$$

So we have

$$1 + 3 + 5 + \dots + (2n-1) = n^2$$

$$1 + 2 + 3 + \dots + n = \frac{1}{2}n^2 + \frac{1}{2}n.$$

Is there any pattern here which might help with our present problem?

These two formulas have one common feature: both are quadratic polynomials in n . Might not the formula we want here also be a polynomial? It seems unlikely that a quadratic polynomial could do the job in this more complicated problem, but how about one of higher degree? Let's try a cubic: assume that there is a formula,

$$1^2 + 2^2 + \dots + n^2 = an^3 + bn^2 + cn + d,$$

where a , b , c , and d are numbers yet to be determined. Substituting $n = 1, 2, 3$, and 4 successively in this formula, we get

$$1^2 = a + b + c + d$$

$$1^2 + 2^2 = 8a + 4b + 2c + d$$

$$1^2 + 2^2 + 3^2 = 27a + 9b + 3c + d$$

$$1^2 + 2^2 + 3^2 + 4^2 = 64a + 16b + 4c + d$$

Solving, we find

$$a = \frac{1}{3}, \quad b = \frac{1}{2}, \quad c = \frac{1}{6}, \quad d = 0.$$

We therefore conjecture that

$$\begin{aligned} 1^2 + 2^2 + \dots + n^2 &= \frac{1}{3} n^3 + \frac{1}{2} n^2 + \frac{1}{6} n \\ &= \frac{1}{6} n(n+1)(2n+1). \end{aligned}$$

This then is our assertion A_n ; now let us prove it.

We have A_k :

$$1^2 + 2^2 + \dots + k^2 = \frac{1}{6} k(k+1)(2k+1).$$

Add $(k+1)^2$ to both sides, factor, and simplify:

$$\begin{aligned} 1^2 + 2^2 + \dots + k^2 + (k+1)^2 &= \frac{1}{6} k(k+1)(2k+1) + (k+1)^2 \\ &= (k+1) \left[\frac{1}{6} k(2k+1) + (k+1) \right] \\ &= \frac{1}{6} (k+1)(k+2)(2k+3), \end{aligned}$$

and this last equation is just A_{k+1} , which is therefore true if A_k is true. Moreover, A_1 , which states

$$1^2 = \frac{1}{6} (1)(2)(3),$$

is true; and A_n is therefore true for each positive integer n .

There is another formulation of the principle of mathematical induction which is extremely useful. This form involves the assumption in the sequential step that every assertion up to a certain point is true, rather than just

the one assertion immediately preceding. Specifically, we have the following (Second Principle of Mathematical Induction): Again let A_1, A_2, A_3, \dots be a sequence of assertions, and let H be the hypothesis that all of these are true. The hypothesis H will be accepted as proved if

1. There is a general proof to show that if every preceding assertion A_1, A_2, \dots, A_k is true, then the next assertion A_{k+1} is true.
2. There is a special proof to show that A_1 is true.

It is not hard to show that either one of the two principles of mathematical induction can be derived from the other. The demonstration of this is left as an exercise.

The value of this second principle of mathematical induction is that it permits the treatment of many problems which would be quite difficult to handle directly on the basis of the first principle. Such problems usually present a more complicated appearance than the kind which yield directly to an attack by the first principle.

Example A3-1h. Every nonempty set S of natural numbers (whether finite or infinite) contains a least element.

Proof. The induction is based on the fact that S contains some natural number. The assertion A_k is that if k is in S , then S contains a least element.

Initial Step: The assertion A_1 is that if S contains 1, then it contains a least number. This is certainly true, since 1 is the smallest natural number and so is smaller than any other member of S .

Sequential Step: We assume A_n is true for all natural numbers up to and including k . Now let S be a set containing $k+1$. There are two possibilities:

1. S contains a natural number p less than $k+1$. In that case p is less than or equal to k . It follows that S contains a least element.
2. S contains no natural number less than $k+1$. In that case $k+1$ is least.

This example is valuable because it is a third principle of mathematical induction equivalent to the other two, although not an obvious one to be sure. An amusing example of a "proof" by this principle is given by Beckenbach in the American Mathematical Monthly, Vol. 52; 1945.

THEOREM. Every natural number is interesting.

Argument. Consider the set S of all uninteresting natural numbers. This set contains a least element. What an interesting number, the smallest in the set of uninteresting numbers! So S contains an interesting number after all. (Contradiction.)

The trouble with this "proof" of course is that we have no definition of "interesting"; one man's interest is another man's boredom.

One of the most important uses of mathematical induction is in definition by recursion, that is, in defining a sequence of things as follows: a definition is given for the initial object of the sequence, and a rule is supplied so that if any term is known the rule provides a definition for the succeeding one.

For example, we could have defined a^n ($a \neq 0$) recursively in the following way:

Initial Step: $a^0 = 1$.

Sequential Step: $a^{k+1} = a \cdot a^k$ ($k = 0, 1, 2, 3, \dots$).

Here is another useful definition by recursion: Let $n!$ denote the product of the first n positive integers. We can define $n!$ recursively as follows:

Initial Step: $1! = 1$.

Sequential Step: $(k+1)! = (k+1)(k!)$ ($k = 1, 2, 3, \dots$).

Such definitions are convenient in proofs by mathematical induction. Here is an example which involves the two definitions we have just given.

Example A3-1i. For all positive integral values n , $2^{n-1} \leq n!$. The proof by mathematical induction is direct. We have the following steps.

Initial Step: $2^0 = 1 \leq 1! = 1$.

Sequential Step: Assuming that the assertion is true at the k -th step, we seek to prove it for the $(k+1)$ -th step. By definition, we have

$$(k+1)! = (k+1)(k!).$$

From the hypothesis, $k! \geq 2^{k-1}$, and consequently,

$$(k+1)! = (k+1)(k!) \geq (k+1)2^{k-1} \geq 2 \cdot 2^{k-1} = 2^k$$

since $k \geq 1$ (k is a positive integer). We conclude that $(k+1)! \geq 2^k$. The proof is complete.

Before we conclude these remarks on mathematical induction, a word of caution. For a complete proof by mathematical induction it is important to show the truth of both the initial step and the sequential step of the induction principle being used. There are many examples of mathematical induction gone haywire because one of these steps fails. Here are two examples.

Example A3-1j.

Assertion: All natural numbers are even.

Argument: For the proof we utilize the second principle of mathematical induction and take for A_k the assertion that all natural numbers less than or equal to k are even. Now consider the natural number $k+1$. Let i be any natural number with $i \leq k$. The number j such that $i+j = k+1$ can easily be shown to be a natural number with $j \leq k$. But if $i \leq k$ and $j \leq k$, both i and j are even; and hence $k+1 = i+j$, the sum of two even numbers, and must itself be even!

Find the hole in this argument.

Example A3-1k.

Assertion: All girls are the same.*

Argument: Given girls designated by a and b , let $a = b$ mean that a and b are the same. Consider any set S_1 containing just one girl. Clearly, if a and b denote girls in S_1 , then $a = b$. Now suppose it is true for any set of k girls that they are all the same. Let S_{k+1} be a set containing $k+1$ girls $g_1, g_2, \dots, g_k, g_{k+1}$. By hypothesis the k girls, g_1, g_2, \dots, g_k , are all the same, but by the same argument so are the k girls $g_2, g_3, \dots, g_k, g_{k+1}$. It follows that $g_1 = g_2 = \dots = g_k = g_{k+1}$. We conclude that all girls of a set containing any positive integral number of them are the same. Since there is only a positive integral number of girls in the whole world, the assertion is proved.

Find the hole in this argument.

*We are not trying to express an overly blasé attitude about girls. The original of this example (attributed to the famous logician Tarski) had it that all positive integers are the same; however, isn't it more interesting to write about girls?

Exercises A3-1

1. Prove by mathematical induction that $1 + 2 + 3 + \dots + n = \frac{1}{2}n(n+1)$.
2. By mathematical induction prove the familiar result, giving the sum of an arithmetic progression to n terms:

$$a + (a + d) + (a + 2d) + \dots + (a + (n-1)d) = \frac{n}{2} [2a + (n-1)d]$$

3. By mathematical induction prove the familiar result, giving the sum of a geometric progression to n terms:

$$a + ar + ar^2 + \dots + ar^{n-1} = \frac{a(r^n - 1)}{r - 1}$$

Prove the following four statements by mathematical induction.

4. $1^2 + 3^2 + 5^2 + \dots + (2n-1)^2 = \frac{1}{3}(4n^3 - n)$.
5. $2n \leq 2^n$.
6. If $p > -1$, then, for every positive integer n , $(1+p)^n \geq 1 + np$.
7. $1 + 2 \cdot 2 + 3 \cdot 2^2 + \dots + n \cdot 2^{n-1} = 1 + (n-1)2^n$.

Prove the following by the second principle of mathematical induction.

8. For all natural numbers n , the number $n+1$ either is a prime or can be factored into primes.
9. For each natural number n greater than one, let U_n be a real number with the property that for at least one pair of natural numbers p, q with $p+q=n$, $U_n = U_p + U_q$.
When $n=1$, we define $U_1 = a$ where a is some given real number.
Prove that $U_n = na$ for all n .
10. Attempt to prove 8 and 9 from the first principle to see what difficulties arise.

In the next three problems, first discover a formula for the sum, and then prove by mathematical induction that you are correct.

$$11. \frac{1}{1 \cdot 2} + \frac{1}{2 \cdot 3} + \frac{1}{3 \cdot 4} + \dots + \frac{1}{n(n+1)}$$

12. $1^3 + 2^3 + 3^3 + \dots + n^3$. (Hint: Compare the sums you get here with Examples A3-1a and A3-1g in the text, or, alternatively, assume that the required result is a polynomial of degree 4.)
13. $1 \cdot 2 + 2 \cdot 3 + 3 \cdot 4 + \dots + n(n+1)$. (Hint: Compare this with Example A3-1g in the text.)
14. Prove for all positive integers n ,

$$\left(1 + \frac{3}{1}\right)\left(1 + \frac{5}{4}\right)\left(1 + \frac{7}{9}\right) \dots \left(1 + \frac{2n+1}{n^2}\right) = (n+1)^2.$$

15. Prove that $(1+x)(1+x^2)(1+x^4) \dots (1+x^{2^n}) = \frac{1-x^{2^{n+1}}}{1-x}$.
16. Prove that $n(n^2+5)$ is divisible by 6 for all integral n .
17. Any infinite straight line separates the plane into two parts; two intersecting straight lines separate the plane into four parts; and three non-concurrent lines, of which no two are parallel, separate the plane into seven parts. Determine the number of parts into which the plane is separated by n straight lines of which no three meet in a single common point and no two are parallel; then prove your result. Can you obtain a more general result when parallelism is permitted? If concurrence is permitted? If both are permitted?
18. Consider the sequence of fractions

$$\frac{1}{1}, \frac{3}{2}, \frac{7}{5}, \frac{17}{12}, \dots, \frac{p_n}{q_n}, \dots$$

where each fraction is obtained from the preceding by the rule

$$p_n = p_{n-1} + 2q_{n-1}$$

$$q_n = p_{n-1} + q_{n-1}.$$

Show that for n sufficiently large, the difference between $\frac{p_n}{q_n}$ and $\sqrt{2}$ can be made as small as desired. Show also that the approximation to $\sqrt{2}$ is improved at each successive stage of the sequence and that the error alternates in sign. Prove also that p_n and q_n are relatively prime, that is, the fraction $\frac{p_n}{q_n}$ is in lowest terms.

19. Let p be any polynomial of degree m . Let $q(n)$ denote the sum

$$(1) \quad q(n) = p(1) + p(2) + p(3) + \dots + p(n).$$

Prove that there is a polynomial q of degree $m + 1$ satisfying (1).

20. Let the function $f(n)$ be defined recursively as follows:

Initial Step: $f(1) = 3$.

Sequential Step: $f(n + 1) = 3^{f(n)}$.

In particular, we have $f(3) = 3^{3^3} = 3^{27}$, etc.

Similarly, $g(n)$ is defined by

Initial Step: $g(1) = 9$.

Sequential Step: $g(n + 1) = 9^{g(n)}$.

Find the minimum value m for each n such that $f(m) \geq g(n)$.

21. Prove for all natural numbers n , that $\frac{(1 + \sqrt{5})^n - (1 - \sqrt{5})^n}{2^n \sqrt{5}}$ is an integer. (Hint: Try to express $x^n - y^n$ in terms of $x^{n-1} - y^{n-1}$, $x^{n-2} - y^{n-2}$, etc.)

A3-2. Sums and Sum Notation.(1) Sum Notation.

In Section 1-1, in the preceding section, and many other places we make frequent use of extended sums in which the terms exhibit a repetitive structure. For example, in Section 1-1 the area of a standard region involves the sum

$$(1) \quad 1 \cdot 1 + 2 \cdot 3 + 3 \cdot 5 + \dots + n(2n - 1).$$

We adopt a concise notation which indicates the repetition instead of spelling it out. In this notation the sum (1) is written

$$\sum_{k=1}^n k(2k - 1).$$

This symbol means, "the sum of all terms of the form $k(2k - 1)$ where k takes on the integer values from 1 to n inclusive." The Greek capital " Σ " (sigma) corresponds to the Roman "S" and is intended to suggest the word "sum."

The notation can be used more generally to express the sum of any quantities ϕ_k where k takes on consecutive integral values; we may begin with any integer m and end with any integer n where $n \geq m$. Thus

$$\sum_{k=m}^n \phi_k = \phi_m + \phi_{m+1} + \phi_{m+2} + \dots + \phi_n.$$

(Note the trivial special case, $n = m$, a "sum" of one term: $\sum_{k=m}^n \phi_k = \phi_m$.)

Example A3-2a. If each of the regions R_k in (1) is a rectangle with height h_k and width w_k , the sum of the areas may be written

$$w_1 h_1 + w_2 h_2 + w_3 h_3 + \dots + w_n h_n = \sum_{k=1}^n w_k h_k.$$

Here are other typical examples:

$$\begin{aligned} \sum_{k=0}^3 \frac{k}{1+k^2} &= \frac{0}{1+0} + \frac{1}{1+1} + \frac{2}{1+4} + \frac{3}{1+9} \\ &= 0 + \frac{1}{2} + \frac{2}{5} + \frac{3}{10} \\ &= \frac{6}{5}, \end{aligned}$$

$$\sum_{j=2}^5 (j+3) = 5 + 6 + 7 + 8 = 26.$$

A linear combination of n functions:

$$\sum_{j=1}^n a_j f_j(x) = a_1 f_1(x) + a_2 f_2(x) + \dots + a_n f_n(x).$$

A polynomial of degree no greater than m :

$$\sum_{i=0}^m c_i x^i = c_0 + c_1 x + c_2 x^2 + \dots + c_m x^m.$$

Example A3-2b. A simple but important sum is $\sum_{j=1}^n c$, where c is a constant, that is, a quantity independent of the index j of summation. The quantity $\sum_{j=1}^n c$ is the sum of n terms each of which is c ; it therefore has the value nc .

In any summation the values of the terms and the total are not affected by the choice of the index letter; thus

$$\sum_{k=m}^n \phi_k = \sum_{j=m}^n \phi_j;$$

We are free to choose the index letter and its initial value to suit our own convenience.

Example A3-2c.

$$(a) \quad \sum_{j=0}^2 a_j = a_0 + a_1 + a_2 = \sum_{p=1}^3 a_{p-1} = \sum_{n=0}^2 a_{2-n}.$$

$$(b) \quad \sum_{i=0}^n a_i^{n-i} = a_0^n + a_1^{n-1} + \dots + a_n^0 = \sum_{j=0}^n (a_{n-j})^j$$

Summation is a linear process; the proof is left as the first exercise below.

Exercises A3-2a

1. Prove

$$\sum_{k=1}^n (\alpha f_k + \beta g_k) = \alpha \sum_{k=1}^n f_k + \beta \sum_{k=1}^n g_k$$

2. Write each of the following sums in expanded form and evaluate:

$$(a) \sum_{k=1}^5 2k$$

$$(d) \sum_{m=2}^5 m(m-1)(m-2)$$

$$(b) \sum_{j=5}^{10} j^2$$

$$(e) \sum_{i=0}^{10} 2^i$$

$$(c) \sum_{r=-1}^3 (r^2 + r - 12)$$

$$(f) \sum_{r=0}^4 \frac{4!}{r!(4-r)!}$$

3. Which of the following statements are true and which are false? Justify your conclusions.

$$(a) \sum_{j=3}^{10} 4 = 7 \cdot 4 = 28$$

$$(b) \sum_{j=m}^n 4 = 4((n-m)+1)$$

$$(c) \sum_{k=1}^{10} k^2 = 10 \sum_{k=1}^9 k^2$$

$$(d) \sum_{k=1}^{1000} k^2 = 5 + \sum_{k=3}^{1000} k^2$$

$$(e) \sum_{k=1}^n k^3 = n^3 + \sum_{j=2}^n (j-1)^3$$

$$(f) \sum_{m=1}^{10} k^2 = \left(\sum_{m=1}^{10} k \right)^2$$

$$(g) \sum_{m=1}^{10} k^3 = \left(\sum_{m=1}^{10} k \right)^2$$

$$(h) \sum_{i=0}^n i(i-1)(n-i) = \sum_{i=2}^{n-1} i(i-1)(n-i)$$

$$(i) \sum_{k=0}^m f(a_{m-k}) = \sum_{k=0}^m f(a_k)$$

$$(j) n \sum_{k=0}^n A_k - \sum_{k=0}^n k A_k = \sum_{k=0}^n k A_{n-k}$$

$$(k) \sum_{k=0}^m k^2 (A_k - A_{m-k}) = m^2 \sum_{k=0}^m A_{m-k} - 2m \sum_{k=0}^m k A_{m-k}$$

4. Evaluate $\sum_{k=1}^n f\left(\frac{k}{n}\right) \frac{(b-a)}{n}$ if $f(x) = x^2$, $a = 0$, $b = 1$, and

(a) $n = 2$.

(b) $n = 4$.

(c) $n = 8$.

5. Subdivide the interval $[0,1]$ into n equal parts. In each sub-interval obtain upper and lower bounds for x^2 . Using sigma notation use these upper and lower bounds to obtain expressions for upper and lower estimates of the area under the curve $g = x^2$ on $[0,1]$. If you can evaluate these sums without reading elsewhere, do so.

6. (a) Write out the sum of the first 7 terms of an arithmetic progression with first term a and common difference d . Express the same sum in sigma notation.

(b) In sigma notation, write the expression for the sum of the first n terms of a geometric progression with first term a and common ratio r .

7. (a) Consider a function f defined by

$$f(n) = \sum_{r=1}^n \{(r-1)(r-2)(r-3)(r-4)(r-5) + r\}.$$

Find $f(n)$ for $n = 1, 2, \dots, 5$.

(b) Give an example of a function g (similar to that in (a)) such that

$$g(n) = 1 \quad n = 1, 2, \dots, 10^6,$$

$$g(10^6 + 1) = 0.$$

8. Write each of the following sums in expanded form and evaluate.

(a) $\sum_{n=1}^4 \left\{ \sum_{r=1}^3 r(n-r) \right\}$

(b) $\sum_{n=1}^N \left\{ \sum_{r=1}^R (rn-1) \right\}$

9. The double sum $\sum_{i=0}^m \sum_{j=0}^n F(i,j)$ is a shorthand notation for

$$\sum_{i=0}^m \{ F(i,0) + F(i,1) + \dots + F(i,n) \}$$

or

$$\begin{aligned} & F(0,0) + F(0,1) + \dots + F(0,n) \\ & + F(1,0) + F(1,1) + \dots + F(1,n) \\ & \vdots \\ & + F(m,0) + F(m,1) + \dots + F(m,n) \end{aligned}$$

In particular $\sum_{i=1}^2 \sum_{j=1}^3 i \cdot j = 1 \cdot 1 + 1 \cdot 2 + 1 \cdot 3 + 2 \cdot 1 + 2 \cdot 2$

$+ 2 \cdot 3 = 18$. Evaluate:

(a) $\sum_{i=1}^m \sum_{j=1}^n i \cdot j$

(c) $\sum_{i=1}^m \sum_{j=1}^n \max\{i,j\}$

(b) $\sum_{i=1}^m \sum_{j=1}^n (i+j)$

(d) $\sum_{i=1}^m \sum_{j=1}^n \min\{i,j\}$

10. (a) Show that $\frac{1}{k(k-1)} = \frac{1}{k-1} - \frac{1}{k}$, $k \neq 0, 1$.

(b) Evaluate $\sum_{k=2}^{1000} \frac{1}{k(k-1)}$.

11. If $S(n) = \sum_{i=1}^n f(i)$, determine $f(m)$ in terms of the sum function S .

12. Determine $f(n)$ in the following summation formulae:

$$(a) \quad 1 = \sum_{i=1}^n f(i)$$

$$(e) \quad \cos nx = \sum_{i=1}^n f(i)$$

$$(b) \quad n = \sum_{i=1}^n f(i)$$

$$(f) \quad \sin(an + b) = \sum_{i=1}^n f(i)$$

$$(c) \quad n^2 = \sum_{i=1}^n f(i)$$

$$(g) \quad n! = \sum_{i=1}^n f(i)$$

$$(d) \quad an^2 + bn + c = \sum_{i=1}^n f(i)$$

13. Binomial Theorem: We define $\binom{n}{r} = \frac{n!}{(n-r)!r!}$ where r, n are integers such that $0 \leq r \leq n$. Also $0! = 1$ and $\binom{n}{r} = 0$ if $r > n$. Show that

$$(a) \quad \binom{n}{0} = \binom{n}{n} = 1$$

$$(b) \quad \binom{n}{r} = \binom{n}{n-r}$$

$$\binom{n}{1} = \binom{n}{n-1} = n$$

$$\binom{n}{r} + \binom{n}{r+1} = \binom{n+1}{r+1}$$

(c) Establish the Binomial Theorem

$$(x+y)^n = \sum_{r=0}^n \binom{n}{r} x^{n-r} y^r = x^n + nx^{n-1}y + \dots + nxy^{n-1} + y^n$$

$n = 0, 1, 2, \dots$, by mathematical induction.

14. Using the Binomial Theorem, give the expansions for the following:

$$(a) \quad (x+y)^3$$

$$(c) \quad (2x-3y)^3$$

$$(b) \quad (x-y)^3$$

$$(d) \quad (x-2y)^5$$

15. Evaluate the following sums.

$$(a) \quad \binom{n}{0} + \binom{n}{1} + \dots + \binom{n}{n} = \sum_{r=0}^n \binom{n}{r}$$

$$(b) \quad \binom{n}{0} - \binom{n}{1} + \dots + (-1)^n \binom{n}{n} = \sum_{r=0}^n (-1)^r \binom{n}{r}$$

16. Sum $\sum_{r=0}^n r \binom{n}{r}$ by first showing $\sum_{r=0}^n r \binom{n}{r} = \sum_{r=0}^n (n-r) \binom{n}{r}$ and using 15(a).

17. If $P_n(x)$ denotes a polynomial of degree n such that $P_n(x) = 2^x$ for $x = 0, 1, 2, \dots, n$ find $P_n(n+1)$.

(ii) Summation.

Exercises A3-1, No. 10 illustrates a particularly useful summation technique, i.e., representation as a telescoping sum. It was possible to write

$$\sum_{k=2}^{1000} \frac{1}{k(k-1)} = \frac{1}{2 \cdot 1} + \frac{1}{3 \cdot 2} + \frac{1}{4 \cdot 3} + \dots + \frac{1}{1000 \cdot 999}$$

in the form

$$\sum_{k=2}^{1000} \left(\frac{1}{k-1} - \frac{1}{k} \right) = (1 - \frac{1}{2}) + (\frac{1}{2} - \frac{1}{3}) + \dots + (\frac{1}{998} - \frac{1}{999}) + (\frac{1}{999} - \frac{1}{1000})$$

Each quantity subtracted in one parenthesis is added back in the next, so that the first two terms telescope from a sum of four numbers to a sum of two numbers, the first three terms telescope from a sum of six numbers to a sum of two numbers, etc. Finally, the entire summation telescopes (or collapses) into a sum of two numbers--the first number in the first term and the second number in the last term. Symbolically, a telescoping sum has the form

$$(1) \quad \sum_{k=m}^n \{f(k) - f(k-1)\} = f(n) - f(m-1)$$

In the above example, we have $m = 2$, $n = 1000$, and $f(k) = -\frac{1}{k}$ so that the sum telescopes to $f(1000) - f(1) = -\frac{1}{1000} + 1 = \frac{999}{1000}$.

We now use (1) to establish a short dictionary of summation formulae by considering different functions $f(k)$. Also, we let $m = 1$ without loss of generality. Let $f(k) = k$, then

$$(2) \quad \sum_{k=1}^n \{k - (k-1)\} = \sum_{k=1}^n 1 = n$$

This result is nothing new. Now let $f(k) = k^2$, then

$$\sum_{k=1}^n \{k^2 - (k-1)^2\} = \sum_{k=1}^n (2k-1) = 2 \sum_{k=1}^n k - \sum_{k=1}^n 1 = n^2$$

or, equivalently,

$$(3) \quad \sum_{k=1}^n k = \frac{1}{2} n(n+1).$$

By linearly combining (2) and (3), we obtain the sum of a general arithmetic progression

$$\sum_{k=1}^n (ak + b) = a \left\{ \frac{n(n+1)}{2} \right\} + bn.$$

To obtain the sum $\sum_{k=1}^n k^2$, we let $f(k) = k^3$. Then,

$$\sum_{k=1}^n \{k^3 - (k-1)^3\} = \sum_{k=1}^n (3k^2 + 3k + 1) = n^3,$$

$$3 \sum_{k=1}^n k^2 - 3 \sum_{k=1}^n k + \sum_{k=1}^n 1 = n^3.$$

Using (2) and (3), we obtain

$$\sum_{k=1}^n k^2 = \frac{1}{3} \left\{ n^3 - \frac{3n(n+1)}{2} + n \right\} = \frac{n(n+1)(2n+1)}{6}.$$

We now can establish a sequential method of obtaining sums of the form

$$\sum_{k=1}^n P(k) \text{ whose terms are values } P(k) \text{ of a polynomial function. Because a}$$

polynomial is a linear combination of powers, and summation is a linear process, it is sufficient to give a sequential method for $\sum_{k=1}^n k^r$, r a nonnegative integer.

Choosing $f(k) = k^{r+1}$ in summation, formula (1) gives us

$$\sum_{k=1}^n \{k^{r+1} - (k-1)^{r+1}\} = n^{r+1}.$$

Using the Binomial Theorem, we obtain

$$(4) \quad k^{r+1} - (k-1)^{r+1} = (r+1)k^r + P(k)$$

where $P(k)$ is a polynomial of degree $r - 1$. Thus, the sum $\sum_{k=1}^n k^r$ can be expressed in terms of sums of lower degree. Since we already have the sum for $r = 0, 1$, and 2 , we can repeat the method sequentially to obtain the sum for any r (compare with Exercises A3-1, No. 19).

We can enlarge our summation table by choosing other functional forms $f(k)$, e.g., $\sin(ak + b)$. By (1),

$$(5) \quad \sum_{k=1}^n \{ \sin(ak + b) - \sin(a(k-1) + b) \} = \sin(an + b) - \sin b.$$

Using the identity

$$\sin A - \sin B = 2 \sin \frac{A - B}{2} \cos \frac{A + B}{2},$$

in Equation (5), we obtain

$$(6) \quad \sum_{k=1}^n \cos(ak + b - \frac{a}{2}) = \cos(b + \frac{an}{2}) \frac{\sin \frac{an}{2}}{\sin \frac{a}{2}}.$$

If $b = \frac{a}{2}$, (6) reduces to

$$(7) \quad \sum_{k=1}^n \cos ak = \cos\left((a+1)\frac{n}{2}\right) \frac{\sin \frac{an}{2}}{\sin \frac{a}{2}}.$$

If $b = \frac{a}{2} + \frac{\pi}{2}$, (6) reduces to

$$(8) \quad \sum_{k=1}^n \sin ak = \sin(b + \frac{an}{2}) \frac{\sin \frac{an}{2}}{\sin \frac{a}{2}}.$$

By choosing other functions $f(k)$, we can enlarge our list of summation formulae. We leave this for exercises.

Exercises A3-2b

1. Write the following sums in telescoping form, i.e., in the form

$$\sum_{k=1}^n \{u(k) - u(k-1)\}, \text{ and evaluate.}$$

(a) $\sum_{k=1}^n k(k+1)$

(e) $\sum_{k=1}^n k^3$

(b) $\sum_{k=1}^n k(2k-1)$

(f) $\sum_{k=1}^n \frac{1}{k(k+1)(k+2)}$

(c) $\sum_{k=1}^n 2k(2k+1)$

(g) $\sum_{k=1}^n k \cdot k!$

(d) $\sum_{k=1}^n k(k+1)(k+2)$

(h) $\sum_{k=1}^n r^k$

2. Using $\sum_{k=1}^n \{u(k) - u(k-1)\} = u(n) - u(0)$, establish a short dictionary of summation formulae by considering the following functions u :

(a) $(a + kd)(a + (k+1)d) \dots (a + (k+p)d)$

(b) The reciprocal of (a).

(c) r^k

(d) kr^k

(e) $k^2 r^k$

(f) $k!$

(g) $(k!)^2$

(h) $\arctan k$

(i) $k \sin k$

3. Simplify:

$$\frac{\sin x + \sin 3x + \dots + \sin ((2n-1)x)}{\cos x + \cos 3x + \dots + \cos ((2n-1)x)}$$

4. Another method for summing $\sum P(k)$ (P , a polynomial) can be obtained by using a special case of problem 2a, i.e.,

$$\sum_{k=1}^n \{(k+1)(k)(k-1) \dots (k-r+1) - (k)(k-1)(k-2) \dots (k-r)\} \\ = (n+1)(n)(n-1) \dots (n-r+1),$$

$$\text{or } \sum_{k=1}^n k(k-1) \dots (k-r+1) = \frac{(n+1)(n)(n-1) \dots (n-r+1)}{r+1}.$$

First, we show how to represent any polynomial $P(k)$ of r^{th} degree in the form

$$(i) \quad P(k) = a_0 + a_1 k + \frac{a_2 k(k-1)}{2!} + \dots + \frac{a_r k(k-1) \dots (k-r+1)}{r!}.$$

If $k=0$, then $a_0 = P(0)$; if $k=1$, then $a_1 = P(1) - P(0)$; if $k=2$, then $a_2 = P(2) - 2P(1) + P(0)$. In general, it can be shown that

$$(ii) \quad a_m = P(m) - \binom{m}{1}P(m-1) + \binom{m}{2}P(m-2) - \dots + (-1)^m P(0), \\ m = 0, 1, \dots, r.$$

Since both sides of (i) are polynomials of degree r and (i) is satisfied for $m = 0, 1, \dots, r$, it must be an identity.

$$\text{Now sum } \sum_{k=1}^n P(k).$$

5. Using Prob. 4, find the following sums:

$$(a) \quad \sum_{k=1}^n k^2.$$

$$(b) \quad \sum_{k=1}^n k^3 - \left(\sum_{k=1}^n k \right)^2$$

$$(c) \quad \sum_{k=1}^n k^4.$$

6. (a) Establish Equation (ii) of Number 4.

- (b) Show that a_m is zero for $m > r$.

Appendix 4

FUNCTIONS CONTINUOUS ON AN INTERVAL

A4-1. The Extreme Value Theorem.

The range of a function may include arbitrarily large numbers. For example, the function defined by $f(x) = \frac{1}{x}$ on the domain $\{x : 0 < x \leq 1\}$ has the property that for any positive number $z > 1$ there is at least one value of f , e.g., $f(\frac{1}{z+1}) = z+1$, which is greater than z . This cannot occur for a continuous function whose domain is a closed interval.

THEOREM A4-1. If f is continuous on the closed interval $[a, b]$ then f is bounded above on the interval; that is, there exists a number N for which $f(x) < N$ for all x in $[a, b]$.

Proof. Suppose that the theorem is false. Then $f(x)$ has no upper bound on $[a, b]$. It follows that f cannot be bounded above on both of the "half-intervals" $[a, \frac{a+b}{2}]$ and $[\frac{a+b}{2}, b]$. Let $[a_1, b_1]$ be a half-interval where $f(x)$ lacks an upper bound. The same argument applied again yields a half-interval $[a_2, b_2]$ of $[a_1, b_1]$ where f has no upper bound. Applying the argument recursively, we obtain a nested set of intervals $\{[a_n, b_n]\}$ for which $[a_{n+1}, b_{n+1}]$ is a half-interval of $[a_n, b_n]$ and such that f has no upper bound on any interval in the set. From the Nested Interval Principle (Section A1-5) there exists at least one point ξ common to all the intervals $[a_n, b_n]$.

Now ξ is a point of $[a, b]$ and therefore f is continuous at ξ . Consequently, for any positive ϵ there is a δ -neighborhood of ξ wherein

$$|f(x) - f(\xi)| < \epsilon$$

for those x in the domain of f . In the δ -neighborhood of ξ , then, $f(x)$ is bounded above:

$$f(x) < f(\xi) + \epsilon.$$

But since $b_n - a_n = \frac{b-a}{2^n}$ and ξ lies in $[a_n, b_n]$, it follows that $[a_n, b_n]$ is contained in the δ -neighborhood of ξ for all n satisfying $\frac{b-a}{2^n} < \delta$. This contradicts our prior conclusion that $f(x)$ has no upper bound in $[a_n, b_n]$. Hence our original supposition that f has no upper bound on $[a, b]$ is false.

Corollary. If f is continuous on the closed interval $[a, b]$, then f is bounded below on the interval; that is, there is a number N such that for all x in $[a, b]$, $f(x) > N$.

When f is bounded both above and below on $[a, b]$ then there is a number N such that, for all x in $[a, b]$, $|f(x)| < N$; and we shall simply say that f is bounded.

Most often it is essential and adequate to know whether a function is bounded. If we push our analysis a little further, however, some interesting problems arise. The best upper (lower) bound for a function is, of course, the least upper bound, M (greatest lower bound, m) of the range of the function. Considering the fact that there are function values arbitrarily close to M we might hope and expect that M is actually a function value. That this need not happen can be seen in the following problem.

Is there a function among the functions defined on the closed interval $[0, b]$ with value a at 0 , linear from 0 to p with value 0 at p , and 0 from p to b (see diagram), such that $L(p) = b - p + \sqrt{a^2 + p^2}$,

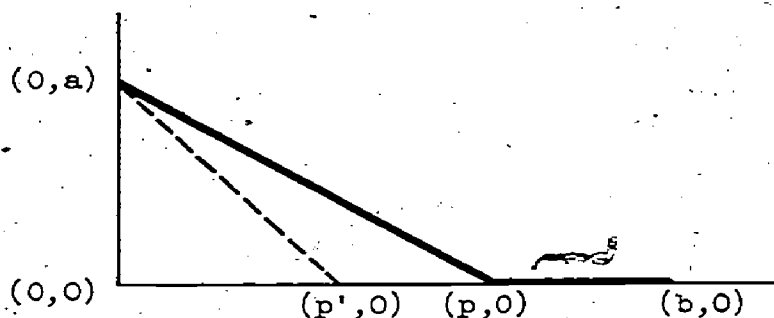


Figure 4-1a

the sum of the length of the line segments that comprise the graph, is as large as possible? If $p' < p$ then the sum of the lengths of the segments

joining $(0,a)$ to $(p',0)$ and $(p',0)$ to $(p,0)$ is greater than the length of the segment joining $(0,a)$ to $(p,0)$. Hence the $L(p') > L(p)$ if $p' < p$. Since p cannot equal 0 (we would not have the graph of a function), the function $L : p \rightarrow L(p)$, though bounded by $a + b$, does not have a maximum value.

The next theorem gives conditions under which this difficulty cannot occur.

THEOREM 3-7b. (Extreme Value Theorem.) If the function f is continuous on the closed interval $[a,b]$, then there are numbers x_m and x_M in $[a,b]$ such that

$$f(x_m) \leq f(x) \leq f(x_M)$$

for all x in $[a,b]$.

Proof 1. Let M be the least upper bound of the set of values $f(x)$ for x in $[a,b]$; then for any number L smaller than M there is a value of f greater than L . In particular, for every positive integer k , there is a point x_k in $[a,b]$ for which $f(x_k) > M - \frac{1}{k}$. Next consider the half-intervals $[a, \frac{a+b}{2}]$ and $[\frac{a+b}{2}, b]$. At least one of the intervals must contain values x_k for infinitely many integers k . (If both only contained x_k for a finite set of integers k , then the whole interval could only contain x_k for a finite set of k 's.) Let $[a_1, b_1]$ be such an interval. Now iterate the process and take for $[a_2, b_2]$ a half-interval of $[a_1, b_1]$ which contains the values x_k for infinitely many integers k .

We define $[a_i, b_i]$ recursively by taking for $[a_{n+1}, b_{n+1}]$ a half-interval of $[a_n, b_n]$ which contains values x_k for infinitely many integers k .

Observe that $\{[a_n, b_n]\}$ is a nested set of closed intervals; that consequently, there is a point ξ common to all of them; and, as in the proof of Theorem A4-1, that any neighborhood of ξ contains $[a_n, b_n]$ provided n is sufficiently large. We now show that $f(\xi) = M$. Since f is continuous, for each $\epsilon > 0$ there exists a $\delta > 0$ such that $|f(x) - f(\xi)| < \epsilon$ for all points x of the domain of f within the δ -neighborhood of ξ . The δ -neighborhood contains some interval $[a_n, b_n]$ of the nested set and therefore contains the values x_k for infinitely many integers k . In particular, the δ -neighborhood must contain some value x_v for which $v > \frac{1}{\epsilon}$, (otherwise

the δ -neighborhood could contain only the values x_k for the finite set of integers satisfying $\left[k \leq \frac{1}{\epsilon}\right]$. Since

$$|f(x_v) - f(\xi)| < \epsilon$$

and

$$f(x_v) > M - \frac{1}{v} > M - \epsilon$$

it follows that

$$M - 2\epsilon < f(\xi) \leq M$$

Since this inequality holds for any positive ϵ we conclude that $f(\xi) = M$; thus the range of the function contains a maximum value.

Proof 2. Since M is the least upper bound of the range of f , for every number M' smaller than M there is a value of f greater than M' . Thus, for every positive integer j there is a point x_j in $[a, b]$ such that $f(x_j) > M - \frac{1}{j}$. Suppose that M is not in the range of f . Then the function $\phi : x \mapsto \frac{1}{M - f(x)}$ is continuous in $[a, b]$ (Theorem 3-6c), hence by the Boundedness Theorem ϕ is bounded. But $\phi(x_j) = \frac{1}{M - f(x_j)} > j$, hence ϕ is not bounded. Contradiction.

Exercises A4-1

1. Is the continuity of f essential to the hypothesis of the boundedness theorem?
2. Can a discontinuous function whose domain is a closed interval be bounded?
3. Do Numbers 1 and 2 amount to the same question?
4. Can a nonconstant function whose domain is the set of real numbers be bounded?
5. Give an example of a monotone function on $[0, 1]$ with exactly n points of discontinuity.

6. Can a monotone function on $[0,1]$ have infinitely many points of discontinuity? Justify your answer.
7. (a) Give an example of a bounded function f defined on $[0,1]$ such that f has no extreme values.
 (b) Repeat (a) with the extra condition that f have an inverse.
8. Give an example of a function f defined in the interval $[0,1]$ such that
 (a) f has neither an upper or lower bound.
 (b) f has a lower bound but no upper bound.
 (c) f achieves the upper and lower bounds an infinite number of times.
 (For this case give a function that is not constant in any interval.)
9. Show that any function f satisfying Number 9(c) cannot be continuous in the entire interval.
10. Give an example of a function f defined on $[0,1]$ such that f takes on every value between 0 and 1 once and only once but is discontinuous for all x .
11. Show that a function which is increasing at every point of an interval (a,b) is an increasing function in (a,b) .
12. A function ϕ is said to be weakly increasing "on the right" at a point a if for all x in a neighborhood $[a, a + \delta]$, $\phi(x) \geq \phi(a)$.
 (a) Show that if ϕ is continuous and weakly increasing on the right of all points in (b,c) then ϕ is weakly increasing in (b,c) .
 (b) Show by a counter-example that (a) doesn't necessarily hold if ϕ is discontinuous.
13. A function has the property that for each point of an interval where it is defined, there is a neighborhood in which the function is bounded. Show that the function is bounded over the whole interval. (This is an example where a local property implies a global one. It is clear that the global property here implies the local one.)
14. Give an example of a function defined everywhere in a closed interval but unbounded in the neighborhood of every point of the interval.
 (Suggestion: See Exercises 3-5, No. 15.)

A4-2. The Intermediate Value Theorem.

The idea of continuity as expressed in the first two paragraphs of Section 3-5 is simple and intuitive. However, in order to attain a precise and workable definition, we abstracted what we thought was the essential property of continuity to give Definition 3-5. How do we know that our definition is appropriate? Is our definition in agreement with our intuitive idea of continuity? Do the functions which satisfy our precise definition have the properties that we want continuous functions to have? Whenever we engage in the process of giving a precise definition of an intuitive idea, we gain evidence for the appropriateness of our definition by "proving the obvious"--obvious in that the property is perceived directly from the intuitive idea, and proven in the sense that it is implied by our precise definition. In Section 3-7 we saw that the Intermediate Value Theorem was obvious. Now we shall prove it.

THEOREM 3-7a. (Intermediate Value Theorem). Let f be continuous on the closed interval $[a, b]$. Let v be any number between $f(a)$ and $f(b)$. Then there is a number u in $[a, b]$ such that $f(u) = v$.

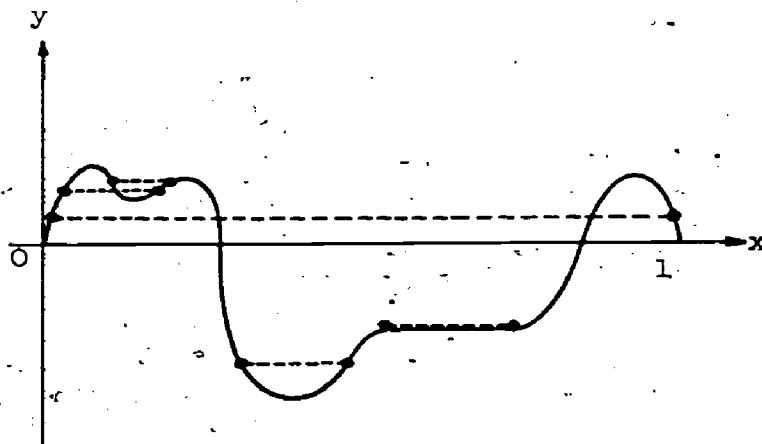
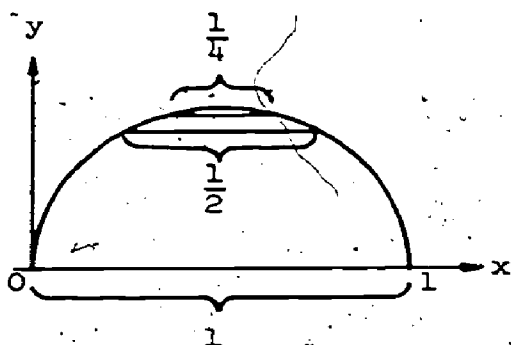
Proof. Suppose $f(a) < v < f(b)$. Let $S = \{x : x \in [a, b] \text{ and } f(x) < v\}$. The set S is not empty since $a \in S$ and it has the upper bound b ; consequently S has a least upper bound u in $[a, b]$. We proceed to show that $f(u) = v$. First, u cannot be an endpoint of $[a, b]$ since, by the continuity of f , $f(x) < v$ in a neighborhood of a and $f(x) > v$ in a neighborhood of b . Next, every neighborhood of u contains points such that $f(x) < v$ (since u is the least upper bound of the set of such points) and points such that $f(x) > v$ (all points to the right of u). It follows that $f(u) = v$, for if $f(u)$ were greater than v then by continuity all points x of some neighborhood of u would satisfy $f(x) > v$ and, similarly, if $f(u) < v$ there would be a neighborhood of u where $f(x) < v$, (cf. Lemma 3-4).

Exercises A4-2

1. Can a discontinuous function have the intermediate value property? Give examples.

2. Let the function f be the derivative of a function g . Prove that f has the intermediate value property.

3. Given the half-circle $y = \sqrt{1 - x^2}$, it can be shown that chords parallel to the x -axis of length $\frac{1}{n}$ exist where n is any positive integer. This result can be generalized to any continuous function taking on the value 0 at 0 and 1. Chords which intersect the curve, or lie entirely outside the curve, or coincide with the curve are permitted. Prove this.



A4-3. A Nowhere Differentiable Continuous Function.

Here we exhibit a Weierstrass Function, that is, an everywhere continuous but nondifferentiable function. The idea is to take a smooth curve and "roughen" it in successive stages of approximation to a curve which is still continuous but everywhere rough.

We shall use polygonal approximations. We roughen the line segment joining two consecutive vertices (a_1, b_1) and (a_2, b_2) , where $a_1 < a_2$, by the following construction. We set $h = \frac{a_2 - a_1}{3}$ and $k = 0.3(b_2 - b_1)$, and replace the original segment by the zigzag line consecutively joining the vertices:

$$(a_1, b_1) ; (a_1 + h, b_2 - k) ; (a_2 - h, b_1 + k) ; (a_2, b_2) .$$

(See Figure A4-3a.)

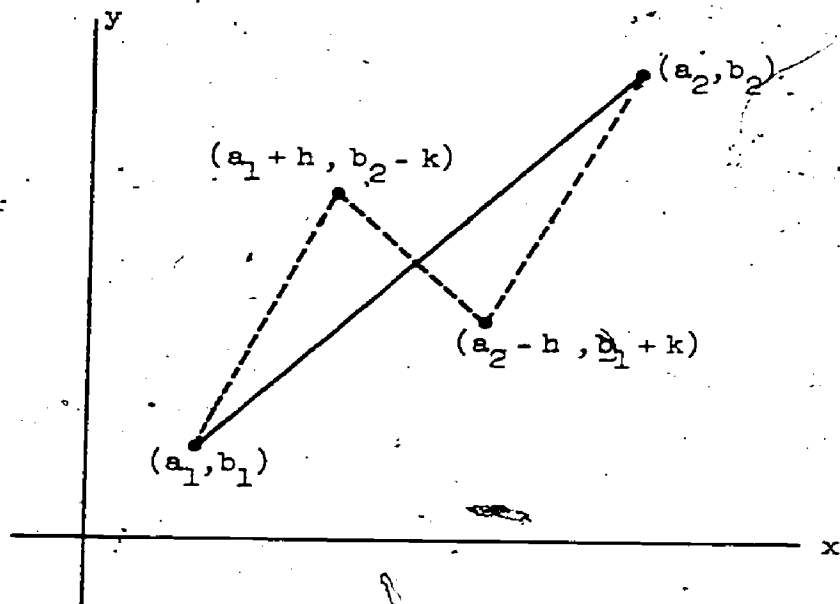


Figure A4-3a

The differences in height of the successive vertices are, in order

$$0.7(b_2 - b_1) , -0.4(b_2 - b_1) , 0.7(b_2 - b_1) .$$

Thus, the absolute difference in height between consecutive vertices is no greater than $0.7(b_2 - b_1)$. In going from one step of approximation to the next, then, we reduce the absolute difference in height between two consecutive vertices by a factor no greater than 0.7 . This is the property which will yield continuity in the limit.

The three consecutive segments of the zigzag line have as their respective slopes, in order,

$$2.1m, -1.2m, 2.1m$$

where $m = \frac{b_2 - b_1}{a_2 - a_1}$ is the slope of the original segment. We see therefore that the steepness or absolute value of the slopes has been multiplied by a factor no less than 1.2 . The increase in steepness of the segments is the property which will yield nondifferentiability in the limit.

Note that as we pass from the segment to the zigzag line, the endpoints and midpoint of the original segment remain fixed.

For the construction of the Weierstrass function f we begin with the segment of the line $y = x$ on the interval $0 \leq x \leq 1$. The next step is to apply the construction above to replace this segment by a zigzag (solid curve in Figure A4-3b). At the next step we replace the segment over each third of the interval by a further application of the construction; all succeeding steps are simply iterations of this procedure. The first five iterations are shown in Figures A4-3b - f.

First we observe that at the n -th step all vertices of the polygonal curve, i.e., the points with abscissas

$$\frac{j}{3^n} \quad (j = 0, 1, 2, \dots, 3^n),$$

remain vertices at the next step. It follows that these points lie on all the approximating polygons from the n -th step on, and hence are points of the prospective limit curve. Therefore, if we let f denote the Weierstrass function, we observe that this construction defines f unambiguously for all the ternary fractions, i.e., for all values $x = \frac{j}{3^n}$, for $n = 1, 2, 3, \dots$, and $j = 0, 1, 2, \dots, 3^n$. We shall use this result to define f unambiguously for all real values as the continuous completion of the function defined only for ternary values.

A4-3

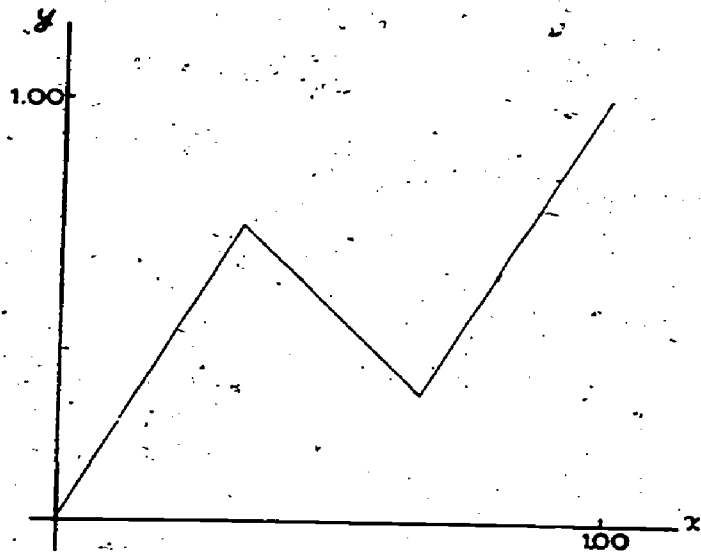


Figure A4-3b

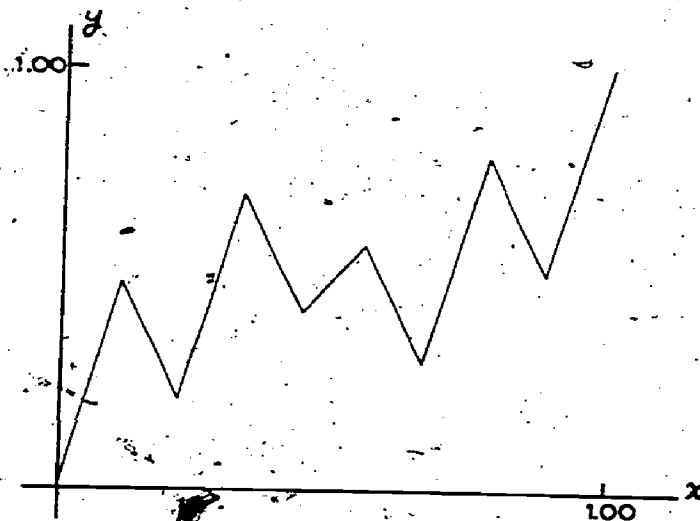


Figure A4-3c

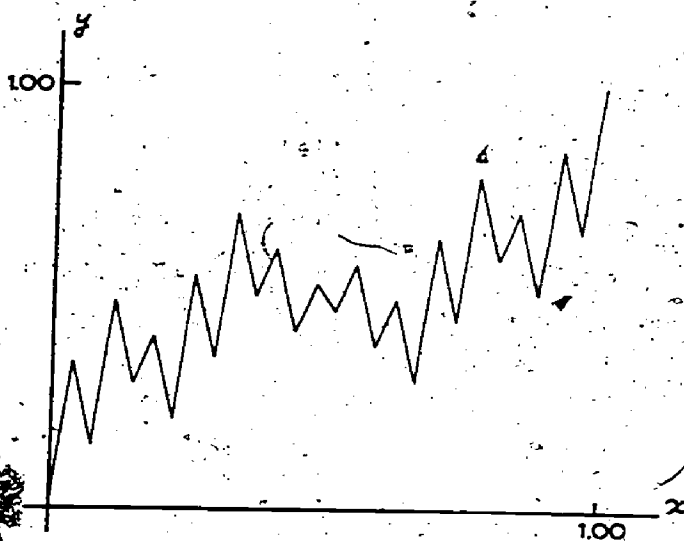


Figure A4-3d

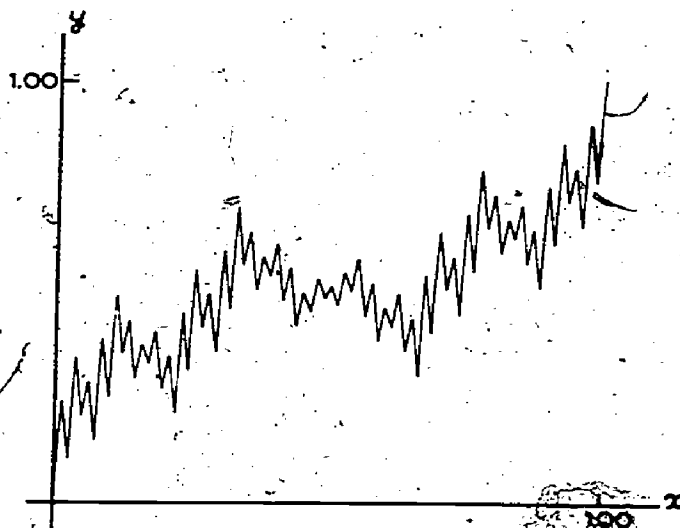


Figure A4-3e

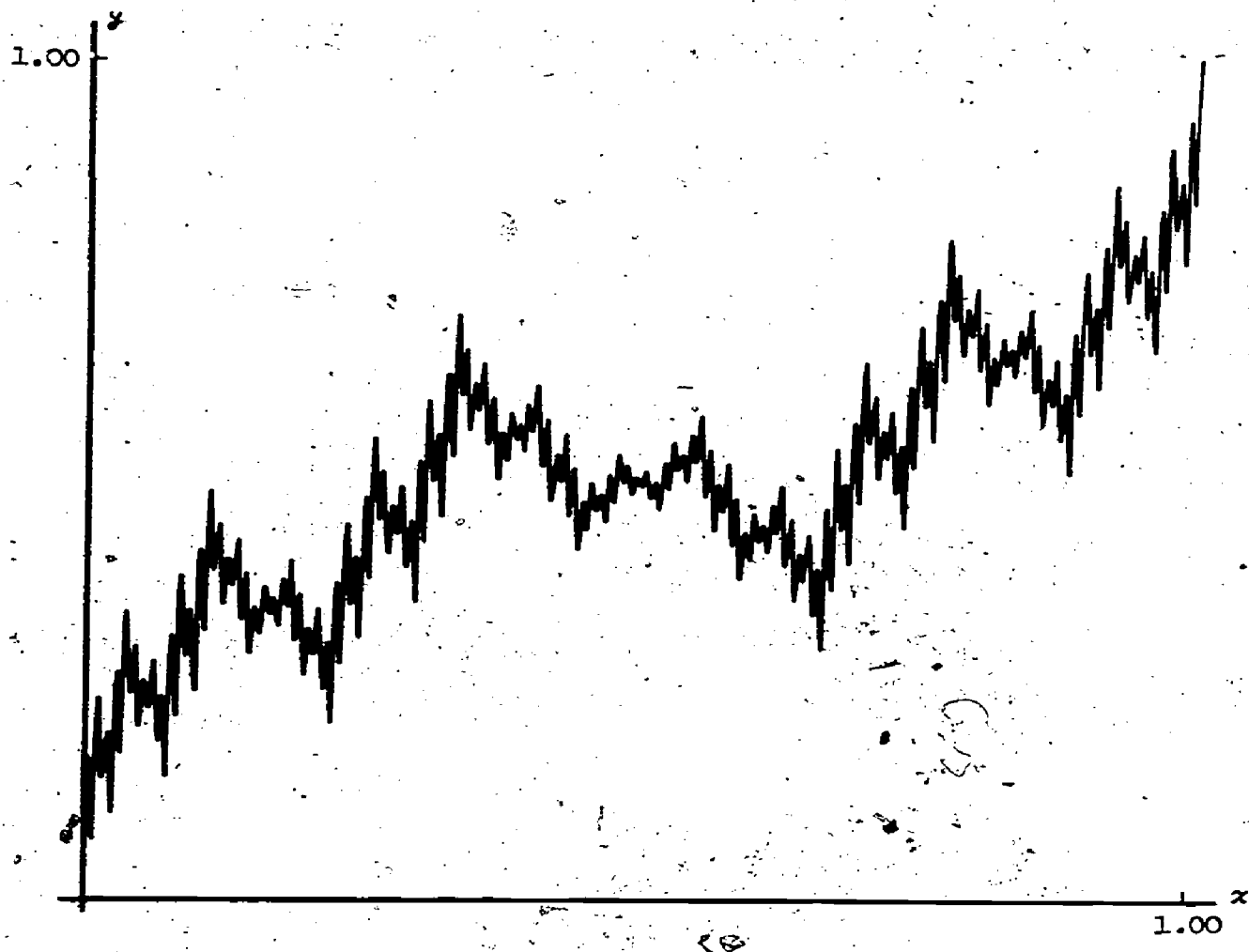


Figure A4-3f

From our discussion of the zigzag construction we observe for two successive vertices x_j and x_{j+1} of the polygon obtained at the n -th iteration that

$$(1) \quad |f(x_{j+1}) - f(x_j)| \leq 0.7^n$$

where $x_j = \frac{j}{3^n}$. This is the basis for the proof of continuity for the limit function. Moreover, for the slope of the chord joining $(x_j, f(x_j))$ to $(x_{j+1}, f(x_{j+1}))$, we have

$$(2) \quad \left| \frac{f(x_{j+1}) - f(x_j)}{x_{j+1} - x_j} \right| \geq (1.2)^n$$

It is clear from (2) that in every subinterval, no matter how small, we can find chords as steep as we please. We shall see that (2) implies that no derivative exists at any point.

Although we have fixed the values of the Weierstrass function at the ternary points we have yet to define $f(x)$ when x is not a ternary point. For each n let j_n be chosen so that

$$a_n \leq x < b_n$$

where $a_n = \frac{j_n}{3^n}$. Observe that $a_{n+1} \geq a_n$ and $b_{n+1} \leq b_n$, and that $f(a_{n+1})$ and $f(b_{n+1})$ both lie between $f(a_n)$ and $f(b_n)$. Thus the set of closed intervals $\{I_n\}$ where I_n has $f(a_n)$ and $f(b_n)$ as endpoints is nested in the sense of Section A1-5. The one real number common to all the intervals I_n is taken as $f(x)$. From this way of defining $f(x)$ it is easy to prove continuity; for $\epsilon > 0$ we choose n so large that $2(0.7)^n < \epsilon$ and take $\delta = 4^{-n}$. If $|\xi - x| < \delta$ and $\xi < x$, then by (1)

$$|f(\xi) - f(x)| \leq |f(\xi) - f(a_n)| + |f(a_n) - f(x)|$$

$$\leq 2(0.7)^n < \epsilon.$$

Similarly, if $|\xi - x| < \delta$ and $\xi > x$,

$$|f(\xi) - f(x)| \leq |f(\xi) - f(b_n)| + |f(b_n) - f(x)| < \epsilon.$$

Now we prove that the function f cannot have a derivative anywhere. For a ternary point $\alpha = \frac{p}{3^q}$, in particular, we can find a ternary point in any neighborhood of α such that the chord joining $(\alpha, f(\alpha))$ to $(\beta, f(\beta))$ has arbitrarily large slope. For this purpose we take $(\alpha, f(\alpha))$ as a vertex on the n -th iterated polygon ($n \geq q$):

$$\alpha = \frac{p3^{n-q}}{3^n}$$

and take

$$\beta = \frac{p3^{n-q} \pm 1}{3^n}$$

Thus we have

$$|\beta - \alpha| = \frac{1}{3^n}$$

and, by (2)

$$\left| \frac{f(\beta) - f(\alpha)}{\beta - \alpha} \right| \geq (1.2)^n.$$

Clearly, by taking n large enough we can obtain a chord with slope as large as we please in any neighborhood of a . For any value r in the interval $0 < r < 1$ which is not a ternary point we can find two successive ternary numbers x_j, x_{j+1} of the subdivision at the n -th step such that

$$x_j < r < x_{j+1}.$$

From (2) we know that $f(x_{j+1}) \neq f(x_j)$, say $f(x_{j+1}) > f(x_j)$ (the argument for $f(x_j) > f(x_{j+1})$ is similar). There are three possibilities (Figure A4-3g):

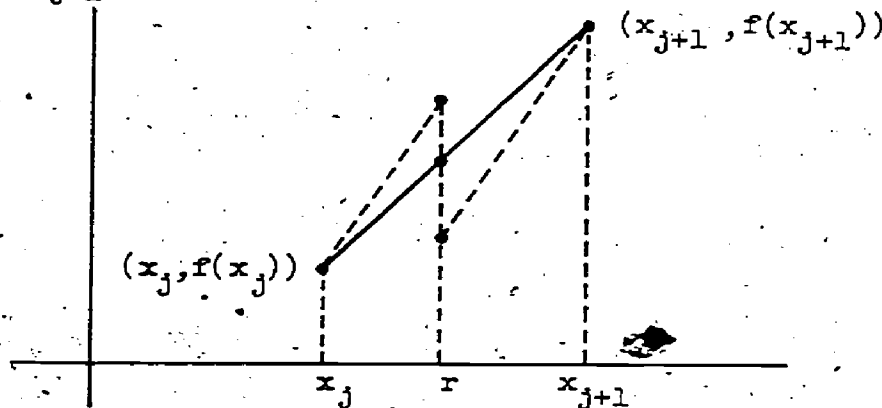


Figure A4-3g

- (1) The point $(r, f(r))$ lies on the chord joining $(x_j, f(x_j))$ to $(x_{j+1}, f(x_{j+1}))$ then

$$\frac{f(x_j) - f(r)}{x_j - r} = \frac{f(x_{j+1}) - f(x_j)}{x_{j+1} - x_j} \geq (1.2)^n$$

- (2) The point lies above the chord; then

$$\frac{f(x_j) - f(r)}{x_j - r} > \frac{f(x_{j+1}) - f(x_j)}{x_{j+1} - x_j} \geq (1.2)^n$$

- (3) The point lies below the chord; then

$$\frac{f(x_{j+1}) - f(r)}{x_{j+1} - r} > \frac{f(x_{j+1}) - f(x_j)}{x_{j+1} - x_j} \geq (1.2)^n$$

Again in any neighborhood of r we can find points for which the slope of the corresponding chords can be made as large as we please.

Exercise A4-3

1. Show that the Weierstrass function is not monotone in any interval.

Appendix 5

IMPLICITLY DEFINED FUNCTIONS AND THEIR DERIVATIVES

In Section 4-8 we gave without proof a method for differentiating implicitly defined functions. Several points remain to be clarified: if a relation between two variables does not express one variable explicitly in terms of another, under what circumstances may it describe a function implicitly? is the implicitly defined function differentiable? if it is differentiable, how may the derivative be found? We consider an example which exhibits most of the difficulties, the equation

$$F(x,y) = (x^2 + y^2)^2 - c^2(x^2 - y^2) = 0$$

whose graph is the lemniscate of Bernoulli (see Figure A5a).

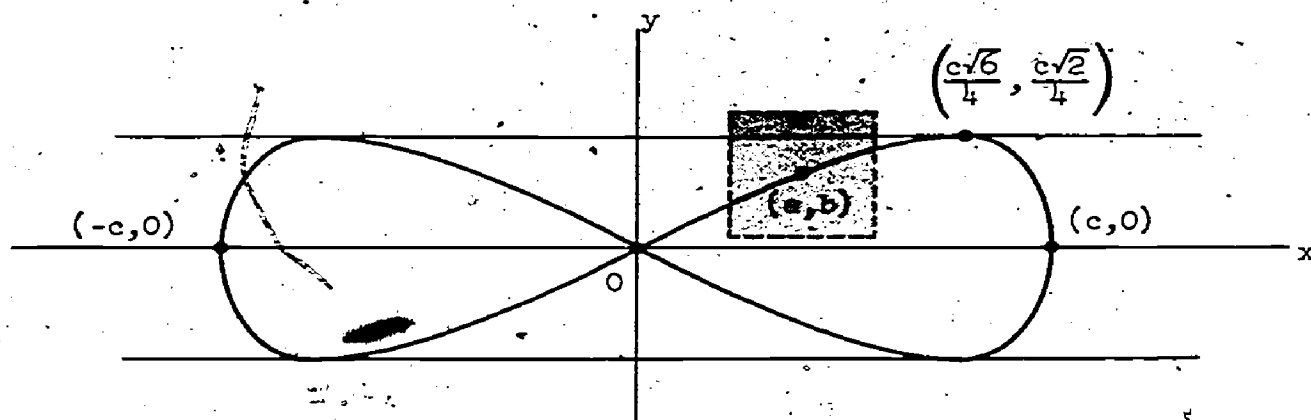


Figure A5a

For each point (a,b) on the graph other than $(0,0)$, $(c,0)$, and $(-c,0)$ there is a neighborhood* of (a,b) wherein no vertical line meets

*The idea of neighborhood of a point of the plane is a natural extension of the idea of neighborhood of a point on the number line. The δ -neighborhood of (a,b) is the square

$$\{(x,y) : |x - a| < \delta \text{ and } |y - b| < \delta\}.$$

In this appendix a δ -neighborhood of a point is to be understood as a linear neighborhood of the point if it is described by one coordinate (point on the number line), as a planar neighborhood of the point if it is described by two coordinates (point of the plane).

the graph more than once (Figure A5a). Within such a neighborhood, then, the graph can be represented by a function $f: x \rightarrow y$, where

$$F(x, f(x)) = (x^2 + [f(x)]^2)^2 - c^2(x^2 - [f(x)]^2) = 0.$$

We cannot do the same thing at $(0,0)$ or at the points $(\pm c, 0)$; at none of these points can the entire graph in a neighborhood be represented by a function.

It is easy to see analytically why the relation $F(x,y) = 0$ must describe a function in the neighborhood of some points of its graph. For example, let us restrict x to the open interval $(0,c)$ and y to be nonnegative. Along a vertical line $x = a$ we introduce the function ϕ_a given by

$$\phi_a(y) = F(a,y) = (a^2 + y^2)^2 - c^2(a^2 - y^2).$$

Now we observe that ϕ_a can have both positive and negative values; namely,

$$\phi_a(a) = F(a,a) = 4a^4 > 0$$

and, since $a < c$,

$$\phi_a(0) = F(a,0) = a^4 - c^2a^2 = a^2(a^2 - c^2) < 0.$$

Since ϕ_a is a polynomial function, hence continuous, we now know by the Intermediate Value Theorem that there is at least one zero of ϕ_a in $(0,a)$; in other terms, the vertical line $x = a$, for $0 < a < c$, must meet the graph $F(x,y) = 0$ in at least one point above the x -axis. In order that the position of the graph $F(x,y) = 0$ in the given region describe a function, each line $x = a$ must meet that part of the graph in no more than one point. This will be the case if ϕ_a has no more than one zero, and, for that, it is sufficient that ϕ_a be strongly monotone. Since, for $y > 0$ and $0 < a < c$,

$$\phi'_a(y) = 2y[2(a^2 + y^2) + c^2] > 0$$

we see that ϕ_a is increasing. Consequently, there is exactly one value $y = b$ for which $\phi_a(b) = F(a,b) = 0$.

In this way we have proved that to each number $x = a$ in $(0,c)$ there corresponds a unique value $b > 0$ such that $F(a,b) = 0$. In other terms, the equation $F(x,y) = 0$, subject to the conditions $x \in (0,c)$ and $y \geq 0$, defines a unique function f such that $F(x, f(x)) = 0$. This same approach to the definition of an implicit function will now be generalized.

THEOREM A5a. (Implicit Function Theorem). Let $F : (x,y) \rightarrow F(x,y)$ be a function defined on a k -neighborhood of $[a,b]$ satisfying the conditions

- (i) $F(a,b) = 0$;
- (ii) For any fixed x , say $x = \xi$ in the k -neighborhood of a , the function $\phi_\xi : y \rightarrow F(\xi,y)$ is increasing and continuous;*
- (iii) For any fixed y , say $y = \eta$ in the k -neighborhood of b , the function $\psi_\eta : x \rightarrow F(x,\eta)$ is continuous.

Under these conditions there exists a unique function $f : x \rightarrow y$ defined on an h -neighborhood ($h \leq k$) of a for which the range of f is contained in the k -neighborhood of b , and

- (1) $f(a) = b$;
- (2) $F(x, f(x)) = 0$,
- (3) f is continuous .

Proof. Since ϕ_a is strongly increasing and $F(a,b) = 0$, it follows that

$$\phi_a(b - \frac{k}{2}) < \phi_a(b) < \phi_a(b + \frac{k}{2})$$

and, since $\phi_a(b) = F(a,b) = 0$, that

$$\phi_a(b - \frac{k}{2}) < 0 < \phi_a(b + \frac{k}{2}) ,$$

or in terms of $F(x,y)$,

$$F(a, b - \frac{k}{2}) < 0 < F(a, b + \frac{k}{2}) .$$

We set $u = b - \frac{k}{2}$ and $v = b + \frac{k}{2}$ and write this inequality in terms of the functions ψ_u and ψ_v :

$$\psi_u(a) < 0 < \psi_v(a) .$$

Since ψ_u and ψ_v are continuous functions of x we conclude from Lemma 3-4 that there is an h -neighborhood of a where both

*If in condition (ii) of the hypothesis, ϕ_ξ is a decreasing function for each ξ , then the theorem applies when $F(x,y)$ is replaced by $-F(x,y)$.

$$\psi_u(x) < 0 \text{ and } \psi_v(x) > 0;$$

that is

$$F(x, b - \frac{k}{2}) < 0 < F(x, b + \frac{k}{2}), \quad \text{for } |x - a| < h;$$

so that

$$\phi_x(b - \frac{k}{2}) < 0 < \phi_x(b + \frac{k}{2}), \quad \text{for } |x - a| < h.$$

Because ϕ_x is a continuous function of y it follows from the Intermediate Value Theorem that for each x satisfying $|x - a| < h$ there is at least one y with $|y - b| < \frac{k}{2}$ for which $\phi_x(y) = 0$; and since ϕ_x is strongly monotone the number y is unique. The function $f : x \mapsto y$ where $\phi_x(y) = 0$ is the unique function satisfying conditions (1) and (2) (Figure A5b).

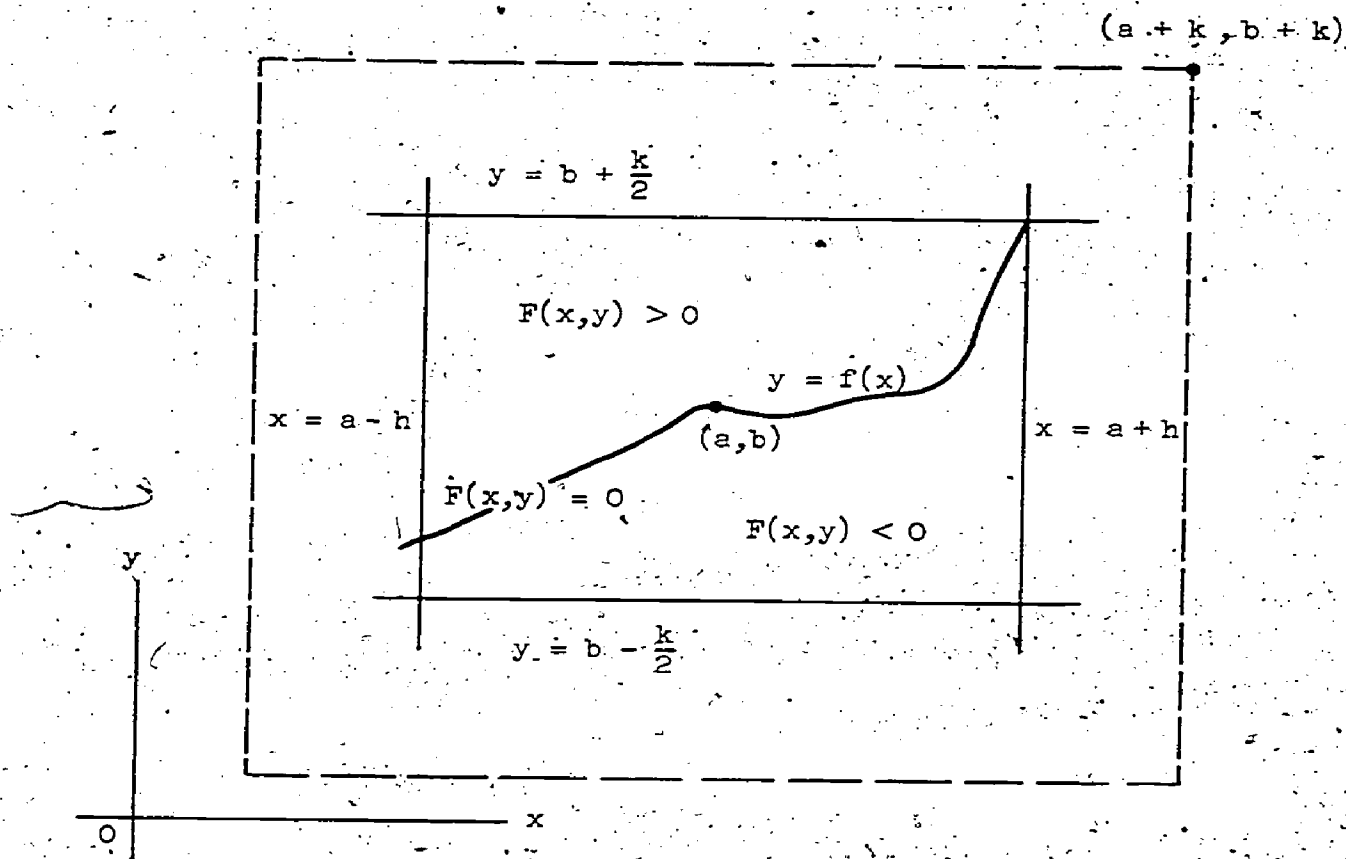


Figure A5b

We have still to show that f is continuous. For any point (α, β) in the k -neighborhood of (a, b) such that $F(\alpha, \beta) = 0$ there exists a k' -neighborhood of (α, β) within the k -neighborhood of (a, b) ; hence, in the k' -neighborhood all the conditions of the theorem and all the conclusions drawn from them above are applicable. In particular, no matter how small k' ,

there exists an h' for which $|f(x) - f(\alpha)| < \frac{k'}{2}$ whenever $|x - \alpha| < h'$, but this is precisely the condition for continuity.

From Section 4-8 we expect that under suitable conditions the implicitly defined function f is not only continuous but that f has a derivative which is easily computed. To prove this, we shall look at the derivative in a way which is useful in many contexts (Section 5-7, also Taylor's Theorem, Chapter 13). We try to approximate a function g locally (near a) by a polynomial. We could, of course, approximate g near a by the constant function $x \rightarrow g(a)$. The next interesting possibility is a linear function and we express $g(x)$ as the sum of a linear term plus an error term; thus

$$(1) \quad g(x) = [\alpha + \beta(x - a)] + e_1(x).$$

If the linear part, $\alpha + \beta(x - a)$, is to be completely significant, then the error term, $e_1(x)$, should be small in comparison with the linear term when $|x - a|$ is small. That is, if we write $e_1(x) = e(x)(x - a)$ we want $\lim_{x \rightarrow a} e(x) = 0$. This is equivalent to the requirement that g be differentiable at a , and that $g(a) = \alpha$ and $g'(a) = \beta$. Thus, the function g is differentiable at a if and only if for some number β ,

$$(2) \quad g(x) = g(a) + \beta(x - a) + e(x)(x - a)$$

and

$$\lim_{x \rightarrow a} e(x) = 0$$

for all x in the domain of g and in some neighborhood of a . This formulation suggests how we may extend many of the properties of linear functions to differentiable functions. Let $g_1(x) = c_1(x - a) + d_1$, $g_2(x) = c_2(x - a) + d_2$, then $g_1'(a) = c_1$ and $g_2'(a) = c_2$. For these linear functions

$$(g_1 + g_2)(x) = (d_1 + d_2) + (c_1 + c_2)(x - a),$$

$$(g_1 \cdot g_2)(x) = d_1 d_2 + (c_1 d_2 + d_1 c_2)(x - a) + c_1 c_2 (x - a)^2,$$

$$g_1 g_2(x) = c_1 (c_2(x - a) + d_2 - a) + d_1 = d_1 + c_1 d_2 - c_1 a + c_1 c_2 (x - a)$$

and we find the appropriate derivatives as the coefficients of $(x - a)$. Moreover, if we think of these linear functions as approximations to other functions, we see that these formulas give the correct coefficients in the general case. Our new perspective provides a method for discovering theorems

as well as the means of proving them.

Thus for a function F of two variables, the requirement of differentiability at the point (a,b) is put in the following form: there exist numbers α and β , functions e_1 and e_2 , such that for all x in some neighborhood of a , all y in some neighborhood of b , F has the representation,

$$(2) \quad F(x,y) = F(a,b) + \alpha(x-a) + \beta(y-b) + e_1(x,y)(x-a) + e_2(x,y)(y-b),$$

where for each ϵ , $|e_1(x,y)|$ and $|e_2(x,y)|$ can be kept less than ϵ by restricting (x,y) to a suitable δ -neighborhood of (a,b) ,

$$|x-a| < \delta \text{ and } |y-b| < \delta.$$

Given this concept of differentiability of F , we may expect that the function f from x to y defined implicitly by the equation $F(x,y) = 0$, will have the same derivative at $x = a$, as the function f_1 from x to y defined implicitly by the equation

$$\alpha(x-a) + \beta(y-b) = 0.$$

For the function f_1 if $\beta \neq 0$, we have

$$y = f_1(x) = \frac{-\alpha(x-a) + \beta b}{\beta}$$

so that $f_1'(a) = -\frac{\alpha}{\beta}$. Thus, we expect that f is differentiable at a , and that $f'(a) = -\frac{\alpha}{\beta}$. This is, in fact, the case since

$$0 = \alpha(x-a) + \beta(f(x)-b) + e_1(x, f(x))(x-a) + e_2(x, f(x))(f(x)-b)$$

we have

$$\frac{f(x) - f(a)}{x - a} = -\frac{\alpha}{\beta} - \frac{e_1(x, f(x))}{\beta} - \frac{e_2(x, f(x))}{\beta} \frac{f(x) - f(a)}{x - a}$$

whence

$$\frac{f(x) - f(a)}{x - a} = \left(-\frac{\alpha}{\beta} - \frac{e_1(x, f(x))}{\beta} \right) \frac{\beta}{\beta + e_2(x, f(x))}$$

Since we already know that f is continuous, e_1 and e_2 are small for x in a neighborhood of a . Thus, $f'(a) = -\frac{\alpha}{\beta}$. Hence we have proved the following theorem.

THEOREM A5b. If $F(x,y)$ satisfies the conditions of Theorem A5a and F is differentiable at (a,b) with $\beta \neq 0$, then the function f defined implicitly by $F(x,y) = 0$ in a neighborhood of (a,b) is differentiable at a , and $f'(a) = -\frac{\alpha}{\beta}$.

Remark. The numbers α , β are (in the notation of the preceding section) $\psi'_b(a)$ and $\phi'_a(b)$ respectively.

Exercises A5

1. Find conditions for which the equations

$$a_1 t + b_1 x + c_1 y + d_1 = 0$$

$$a_2 t + b_2 x + c_2 y + d_2 = 0$$

determine x and y as functions of t .

2. Consider linear approximations to F and G ; then use Number 1 to generalize Theorem A5b to the situation

$$F(t,x,y) = 0 = G(t,x,y).$$

Propose an expression for $D_t x$.

3. Prove that the following equations have unique solutions for y near the points indicated:

(a) $3x^2 + xy + 3y^2 = 7$, $(\sqrt{\frac{7}{3}}, 0)$.

(b) $x^3 - 3xy + y^3 = 0$, $(\frac{3}{2}, \frac{3}{2})$.

(c) $x \cos xy = 0$, $(1, \frac{\pi}{2})$.

4. Find the first and second derivatives of the solutions (a), (b), and (c) of Number 3.
5. Show that there is no unique solution for y in Number 3 (b) near the point $(2^{2/3}, 2^{1/3})$.
6. For possible inverse functions of $f: x \rightarrow y$, show that

$$\frac{d^2 x}{dy^2} = -\frac{\frac{d^2 y}{dx^2}}{(\frac{dy}{dx})^3}$$

INDEX

- acceleration, 409
- antiderivative, 427
- arccos, 144
- arclength, 385
- arcsin, 143
- arctan, 144
- area function
 - additive property, 367
 - order property, 367
- asymptote, 230
 - horizontal, 232, 234
 - oblique (slant), 233, 234
 - vertical, 232, 234
- attenuation equation, 502

- Binomial Theorem, 338
- bounded growth, 512
- bounded set of points, 261
 - greatest lower bound, 266
 - least upper bound, 265
- bounded variation, 649
- braking coefficient, 513
- Buniakowsky-Schwarz inequality, 403

- catenary, 489
- Cauchy's inequality, 253, 403
- chronaxie τ , 509
- composition of functions, 285f, 102
- conic section, 313
 - directrix, 313
 - eccentricity, 313
 - focus, 313
- constrained extreme value problems, 213
- continuity
 - of composite function (Th. 3-6e), 103
 - of differentiable function (Th. 3-6d), 101
 - on the interval, 108
 - intuitive idea, 62
 - of inverse function (Th. 3-6f), 104
 - piecewise, 589
 - of product of continuous functions (Th. 3-6b), 99
 - of quotient of continuous functions (Th. 3-6c), 100
 - of sum of continuous functions (Th. 3-6a), 99
- convex set, 207
- convexity, 206
 - flexed downward, 207, 234
 - flexed upward, 207, 208, 234
- cosine integral, 635
- cover of an interval, 645

- decay coefficient, 499
- decomposition into partial fraction, 563
- decreasing function, 299, 234
 - weakly, 299, 196
- derivative
 - of a^x , 466
 - of $\arccos x$, 147
 - of $\arcsin x$, 147
 - of $\arctan x$, 147
 - of compositions (Chain Rule) (Th. 4-6), 149
 - of $\cos x$, 139
 - of $\cot x$, 139
 - D_x , 117
 - of e^x , 465
 - of $f: x \rightarrow c$, 118
 - of $f: x \rightarrow x$, 118
 - of $f: x \rightarrow x^2$, 118
 - of $f: x \rightarrow \sqrt{x}$, 118
 - of $f: x \rightarrow \frac{1}{x}$, 118
 - of $f: x \rightarrow |x|$, 118
 - of f' , 117
 - of a function at a point, 49
 - of inverse of differentiable function (Th. 4-3), 132
 - of linear combination (Th. 4-2a), 120
 - of $\log x$, 465
 - of polynomial (Th. 4-2c, Cor. 2), 125
 - of polynomial of differentiable function (Th. 4-2c, Cor. 3), 125
 - power rule for positive integers (Th. 4-2c), 125
 - of a product (Th. 4-2b), 122
 - of quotient of differentiable function (Th. 4-2d, Cor. 1), 128
 - of rational function (Th. 4-2d, Cor. 2), 129
 - of reciprocal of differentiable function (Th. 4-2d), 128
 - of right-hand and left-hand, 121
 - of $\sin x$, 139
 - successive higher, 159
 - of $\tan x$, 139
- differential equations, 429
 - e^x (Th. 8-5a), 471
 - $\sin x$, $\cos x$ (Th. 8-5b), 472
- direction angle, 30
- displacement, total, 408
- domain of a function, 269

- e, 461, 480
 - properties of, 477
- ellipse, 313
 - focal chord, 314
 - latus rectum, 314
- energy density, 503
- epsilonics, 67f.
- exponent
 - definition of zero exponent, 446
 - general laws for negative integers, 446
 - general laws for positive integers, 445
 - rational exponents, 447
- exponential function, 447
 - derivative of, 448
 - inverse function, 448
- exponentially damped sinusoid, 607
- Extreme Value Theorem (Th. 3-7b), 109
 - proof, 347
- extremum, 173
 - isolated, 199
 - local, 176, 181, 200
 - on open interval (Lemma 5-2), 178
 - relative, 176
- field, 245
- function
 - absolute value, 95, 274
 - composite, 286
 - even and odd, 276
 - integer part, 57, 275
 - one-to-one, 290
 - periodic, 277
 - signum (sgn), 276, 61, 62
- function definition, 269
 - circular, 303, 137
 - constant, 274
 - explicitly defined, 162
 - identity, 274
 - implicitly defined, 161
 - inverse circular, 143f.
- Fundamental Theorem of calculus, 425
- global properties of f, 169
- graph sketching, 229, 233
- Green's function, 616
- growth coefficient, 497, 513
- half-life, 499
- Heine-Borel Principle, 645
- hyperbola, 313
- hyperbolic functions, 485
 - cosh x, 485
 - derivatives of, 485
 - inverse, 488
 - sinh x, 485
 - tanh x, 485
- hyperbolic sector, 487
- implicit differentiation, 162
- Implicit Function Theorem, 361
- increasing function, 299, 110, 234
 - weakly, 299, 196
- indefinite integral, 427
- initial value, 497
- initial value problem, 430
- integral
 - continuous function, 648
 - definition, 377
 - estimates of, 437
 - existence, 638
 - Existence Theorem (Th. 6-3a), 378
 - geometric properties, 388
 - limit of Tiemann sum, 383, 643
 - of monotone function (Th. 6-3b), 379
- integral operator, 617
- integrals
 - convergent, 582
 - definite, 570
 - definition, 581
 - divergent, 582
 - improper, 578
 - symmetric, 571
- integration, 535
 - of constant times integrable function, 394
 - of linear combination of integrable functions, 393
 - by parts, 554
 - of a polynomial, 633
 - of rational functions, 563
 - special reductions, 573
 - substitution of circular functions, 546
 - Substitution Rule (Th. 10-2), 540
 - of sum of integrable functions, 395
- Intermediate Value Theorem (Th. 3-7a), 109
 - proof, 350
- interval, 259
 - closed, 259, 109
 - interior point of, 259
 - length of, 259
 - midpoint of, 259
 - open, 259, 109
- inverse function, 131, 291f.
- Lagrange rule of variation of parameters, 615
- latent period, 509
- Law of the Mean, 186, 190

lemniscate of Bernoulli, 313 359

limit

of f at a , 58f.

right-hand and left-hand, 90, 578.

$\frac{\sin x}{x}$, 138

limit theorems

constant function (Th. 3-4a), 79

constant multiple of a function
(Th. 3-4b), 79

linear combination of functions
(Th. 3-4c, Cor.), 81

non-negative function (Lem. 3-4
Cor. 2), 84

product of functions (Th. 3-4d), 82

rational function (Th. 3-4e,
Cor. 2), 86

reciprocal of function
(Th. 3-4e), 85

Sandwich Theorem (Th. 3-4f,
Cor. 1), 86

Squeeze Theorem (Th. 3-4f,
Cor. 2), 87

sum of functions (Th. 3-4c), 80

linear approximation to f , 223

linear differential equation of first
order, 590

forcing term, 591

fundamental solution, 594

general solution, 594

initial value problem, 592

nonhomogeneous equation, 595

reduced equation, 591

linear differential equation of
second order, 603

homogeneous equation, 604

superposition principle, 604

local property of a function, 108

logarithms

base e , 461

base 10 (common), 461

derivative, 449

function, 448

as an integral, 452

logistics equation, 513

lower sum over σ , 376

mapping, 270

mathematical induction, 319

first principle, 323

second principle, 327

maximum, local, 177, 181, 198, 205, 234

minimum, local, 177, 181, 198, 205, 234

mean life-time, 499

Mean Value Theorem of integral

calculus, 402

method of equated coefficients, 566

model

for growth, 497

for decay, 499

one function, 299, 196, 415

inverse of strongly monotone
function (Th. A2-4), 300

linear combinations of, 415

piecewise, 415

sectionally, 415

strongly, 299, 178

neighborhood, 260, 58

deleted, 260, 58

of ∞ , 587

radius of, 260

nested interval principle, 265

nonhomogeneous equation, 595

norm of the partition, 379

normal at a point, 226

notation

D_x , 117

Δ (difference), 156

Δ (increment), 149

$\frac{dy}{dx}$, 156

f' , 117

Leibnizian, 156

orthogonal trajectories, 622

parabola, 313

parameter, 44

partition of $[a, b]$, 376

piecewise continuous, 589

piecewise monotone, 415

point of inflection, 230, 234

polar axis, 308

polar coordinates, 308

primitive of f , 427

radioactive decay, 500

radius vector, 308

range of a function, 269

real numbers

algebraic properties of, 245

σ order relations of, 249

rectifiable, 651

recurrence relations, 558

rheobase, 508

Riemann sum, 381

limit of, 383

Rolle's Theorem (Lemma 5-3), 187

scattering coefficient, 503
second derivative, 205
separable differential equation, 621
Separation Axiom, 263
slope, 30
standard region
 lower bound, 370
 upper bound, 370
Stirling's formula, 482
sum notation, 333, 371
summation, 339
superposition principle, 604
supremum, 265
symbol
 $\max \{r_1, r_2, \dots, r_n\}$, 255
 $\min \{r_1, r_2, \dots, r_n\}$, 65, 72, 74,
 83, 255
symmetry, 570

tangent to the curve, 223
tolerance ϵ (error), 32, 63
triangle inequality, 255

upper sum over σ , 377

velocity
 average, 42
 instantaneous, 42
volume of solid of revolution, 405

Wallis's Product for $\frac{\pi}{2}$, 575
Weierstrass function 352, 111